

1 **High-Rate Displacement Monitoring with Low-Cost Multi-Vision Cameras using Global-**
2 **Local Deep Deblurring and Rauch-Tung-Striebel Smoother**

3 Peng “Patrick” Sun¹⁺, Mohammad Vasef²⁺, Lin Chen³

4 ¹ Assistant Professor, Department of Civil, Environmental, and Construction Engineering, University of Central
5 Florida, 12800 Pegasus Drive, Suite 211, Orlando, FL 32816, email: peng.sun@ucf.edu, (corresponding author)

6 ² Research Assistant, Department of Civil, Environmental, and Construction Engineering, University of Central
7 Florida, 12800 Pegasus Drive, Suite 211, Orlando, FL 32816

8 ³ Associate Professor, Department of Bridge Engineering, Tongji University, Shanghai 200092, China

9 ⁺ co-first authors

10 **Abstract:**

11 Measuring structural vibrations help assess dynamic performances of civil structures and
12 infrastructure. Although conventional displacement sensors have been widely adopted, they are
13 contact-based methods which lack scalability. Recently, computer vision (CV) has been applied
14 as a noncontact method to measure displacements. However, the speed of structural vibration
15 (e.g., in shake table tests) can inevitably cause motion blur that imposes challenges in all image-
16 based object/feature detections, especially for normal portable cameras (without high-speed
17 shutters). To address such issue, the study proposed a multi-vision, full-field sensing framework
18 with affordable cameras using a novel global-local detection and deblurring (GLDD) module,
19 which was designed with a generative adversarial network (GAN)-based deblurring model to
20 enhance detection efficiency and accuracy by restoring blemished videos from multiple
21 perspectives. Rauch-Tung-Striebel (RTS) smoother was studied for data fitting using incomplete
22 observations caused due to severe motion-induced blurs. A shake table test was conducted on an
23 aluminum frame with cameras and conventional sensors monitoring the structural vibrations.
24 Fiducial markers were used to track the movement of the key locations on the structure. Results
25 showed that the proposed method is satisfactory to monitor shake table tests when compared to
26 conventional measurements with root-mean-square errors of 0.51-0.95 mm. The proposed
27 deblurring module restored misdetection by 92.1%, 50.6%, and 25.2% for mild-, medium-, and
28 severe-level motion blurs, respectively. Smoother-based data fitting outperformed filter-based one
29 when dealing with highly blemished images. The proposed monitoring system with GLDD and
30 RTS smoother-based data fitting provides a robust measurement solution when dealing with
31 motion blurs.

32 **Keywords:**

33 shake table test, motion blur, computer vision, generative adversarial networks, structural health
34 monitoring.

35
36

37 **Highlights:**

- 38 ▪ proposed a multi-vision displacement measurement approach with global-local detection and
39 deblurring (GLDD) module using GAN-based deep deblurring method to address the motion
40 blur issue
- 41 ▪ developed an automated algorithm for affordable cameras to monitor displacement in shake
42 table tests, including feature detection, global-local image deblurring, multi-video
43 synchronization, and filter/smoothing-based data fitting
- 44 ▪ studied the performances of different data fitting methods on displacement measurements for
45 severe motion blur cases using Kalman filter and Rauch-Tung-Striebel (RTS) smoother
- 46 ▪ provided the guidelines for using the proposed approach and affordable cameras to achieve
47 displacement monitoring in shake table tests

48 **1. Introduction**

49 Monitoring structural responses (e.g., displacement, acceleration, strain) is used to assess the
50 behavior of civil structures. Measured data from experimental tests (e.g., quasistatic test, shake
51 table test) are usually influenced by the characteristics and limitations of the adopted measurement
52 methods (Zona, 2020). Structural responses are commonly measured using wired, contact sensors
53 at desired locations of a structure. Non-contact measurement methods take one step further by
54 avoiding the physical contact between sensor and structures, such as strain sensors using computer
55 vision (CV) techniques (e.g., digital image correlation (del Rey Castillo et al., 2019),
56 photoluminescence techniques (Sun et al., 2019), and laser Doppler effect (Xu et al., 2019). In
57 addition, existing displacement measurement methods include linear variable differential
58 transformer (LVDT), real-time kinematic (RTK) global navigation satellite systems (GNSS)
59 /global positioning system (GPS) sensors (Bezcioglu et al., 2023), terrestrial laser scanner (Kogut
60 & Pilecka, 2020), and double-integration from acceleration (Zheng et al., 2019). However, these
61 displacement measurement methods exhibit specific limitations, such as low-sampling rate (Ma et
62 al., 2022) of RTK-GNSS, limited accuracy in GPS measurements (Rychlicki et al., 2020), high-

63 noise level in terrestrial laser scanner (Muralikrishnan, 2021), (potential) large low-frequency drift
64 using double integration of accelerations (Zheng et al., 2019), and deployment cost of laser
65 Doppler-based method (Chu, 2005). In addition, accessibility issues (especially in long bridges
66 and high-rise buildings), cost escalation for up-scale measurement, range constraint, and generally
67 the requirement for a stable installation platform are the complexities to consider when utilizing
68 LVDTs.

69 Vision-based methods were studied to obtain displacement measurement and overcome some
70 of the limitations. In recent years, the technological progress in computing power, computer vision
71 algorithms (Sun et al., 2022), and high-speed cameras (Zhang et al., 2016) attracted more attentions
72 on the direct measurement methods (Greenbaum et al., 2016) and further applications on the
73 measurements [e.g., system identification (Yang et al., 2019), finite element model updating (Dong
74 et al., 2020), damage detection (Guo et al., 2019)] of vision-based methods. Vision-based
75 applications in shake table tests started from early 2010's by adopting early-stage feature detection
76 algorithms, (large-size) primitive artificial tags, and localization methods to measure structural
77 displacements (Choi et al., 2011). Structural vibration of full-scale civil infrastructures or large
78 scaled models are usually neither in high speed nor in high frequency [e.g., frequency range for
79 most civil infrastructure is well below 70 Hz (Zona, 2020) or even much lower as several Hz].
80 Therefore, most of the time a portable camera with a low frequency capacity [e.g., 30 frame-per-
81 second (fps)] is sufficient and the blur due to structural vibration will not be a serious issue for
82 displacement monitoring. Most of the current studies focused on the further structural health
83 monitoring (SHM) applications of CV-based displacement measurement (e.g., behavior analysis,
84 load estimation, modal identification, model updating, damage detection) (Dong & Catbas, 2021)
85 and much fewer studies focused on solving practical issues in monitoring applications in shake
86 table tests, such as perspective selection (e.g., single-vision, dual-vision), camera location/pose
87 limitation, illumination condition, occlusion, video frame asynchronization (if there is multiple
88 cameras), and motion-induced image blur.

89 Some of these issues can be addressed in a controlled lab environment during a shake table
90 tests, for example, using proficient direct current (DC) lights to provide adequate illumination and
91 using post-synchronization technique to solve asynchronization issue. However, some other
92 challenges remain to be resolved. For example, experimental studies on structural dynamics in
93 particular, can suffer from motion-induced image blurs. Because shake table tests are conducted

94 in lab environments and researchers may use smaller-scale models subjected to more intense
 95 excitations especially when near resonance, making collected video data much more susceptible
 96 to the issue of motion-induced blur. Most of time researchers would adopt pricy, high-speed
 97 cameras (~\$8-30k, 200-2000 fps) to avoid the motion blurs in their CV applications for shake table
 98 without solving the issue. However, even with the most advanced camera with high-speed shutter,
 99 the motion blur issue is still there when dealing with any fast-moving object relative to the camera
 100 shutter. The studies on image deblurring using post-processing techniques for motion-induced
 101 blurs are found to be very rare if there is any. Hence, a remedy solution is in great need to address
 102 the negative impacts from motion blurs for most current shake table users with affordable
 103 measurement setups, such as portable cameras with normal speed (~\$1-2k, 30-60 fps).

104 The objective of this study was to develop a vision-based displacement monitoring framework
 105 for shake table tests that is robust to the motion blur issue. The study proposed a multi-vision
 106 approach with the ability to remediate motion-blur effect using deblurring module and data fitting
 107 module for accurate displacement monitoring. This paper firstly introduced a multi-vision sensing
 108 approach with a global-local detection and deblurring (GLDD) module to reduce the effect of
 109 motion blur and a data fitting module to estimate midsection based on incomplete observation.
 110 Secondly, the study designed a shake table test to evaluate the proposed sensing approach on an
 111 aluminum frame structure with different severity levels of vibration. Further, the discussion based
 112 on the augmented measurements and data fitting was conducted and useful guidelines was
 113 provided as well. In the end, the paper provided a summary for this work as well as its limitation
 114 and future work.

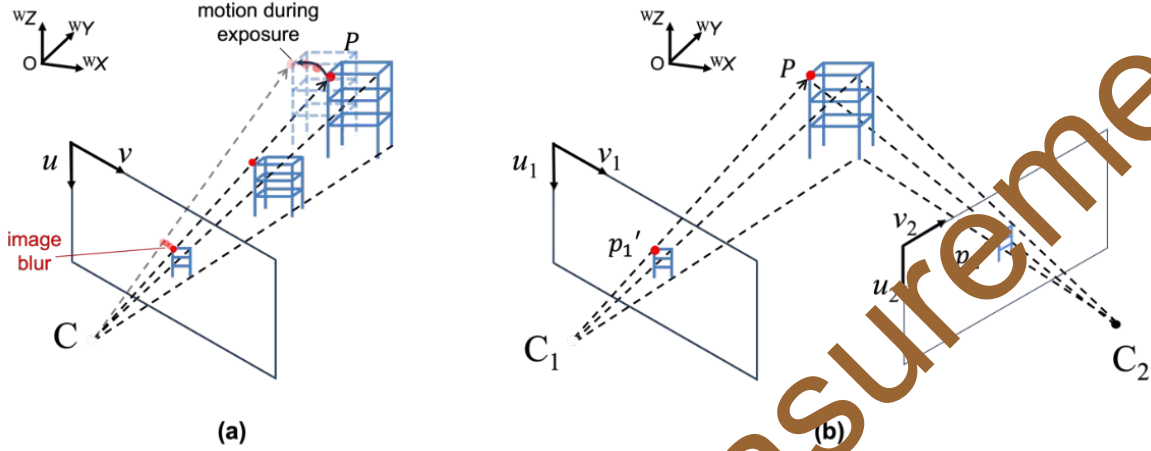
115 **2. Problem Statement**

116 *(1) Motion-Induced Blur in Images*

117 Typical digital image generation includes two steps: image signal acquisition and color
 118 rendering with image signal processor (ISP). In a classical pin-hole camera model (Ma et al., 2004)
 119 (Fig. 2), a 3D point in the World Coordinate System (WCS) is denoted as P and the point is
 120 visualized as a pixel (p') projected onto the camera sensor plane in a 2D Sensor Coordinate System
 121 (SCS). The relationship of P and p' can be expressed in a concise form as:

$${}^s\mathbf{x}_{p'} \sim \mathbf{M}_{\text{aff}} \mathbf{M}_{\text{proj}} [\mathbf{R} | \mathbf{T}] {}^w\mathbf{X}_P \quad (1)$$

122 where M_{aff} and M_{proj} are affine matrix and projection matrix containing intrinsic parameters,
 123 $[\mathbf{R}|\mathbf{T}]$ is the joint rotation-translation matrix (containing extrinsic parameters), ${}^s\mathbf{x}_{p'} =$
 124 $[{}^s x_{p'}, {}^s y_{p'}, 1]^T$ and ${}^w\mathbf{X}_P = [{}^w X_P, {}^w Y_P, {}^w Z_P, 1]^T$ denote the corresponding homogeneous
 125 coordinates, respectively.



126
 127 **Figure 1.** Schematic views of (a) single-perspective setting with motion-induced blur on the sensor
 128 coordinate system (SCS), and (b) dual-perspective setting with the same structure feature observed by both
 129 cameras.

130 Image blurs can result from the relative movement between camera and object/scene and it
 131 can be formulated as the accumulation of photons on camera sensors during the exposure time:

$$I_b(x, y) = \text{ISP} \left(\int_{t_1}^{t_2} f(t, x, y) dt \right) \quad (2)$$

132 where I_b denotes the blurred image, $[t_1, t_2]$ represents the time window for exposure, $f(t, x, y)$
 133 represents the photon response at pixel location (x, y) at time instant t , and $\text{ISP}(\cdot)$ denotes the
 134 image signal processor operator (e.g., white balance, color correction).

135 Motion-induced image blurring (denoted in **Figure 1a**) is influenced by the shutter speeds of
 136 cameras and the relative movement speeds between cameras and recorded objects. Motion blurs
 137 can be categorized into local blurs and global blurs. Global blurs usually occurs with a moving
 138 camera that is usually encountered in the application of robotic vision (Zeng et al., 2020) and
 139 simultaneous localization and mapping (SLAM) (Gao & Zhang, 2021). While local blurring results
 140 from moving objects in static backgrounds. Local blur issue occurs in the CV application for shake
 141 table tests where the table base and the mounted dynamic structures are the foreground in motion
 142 and cameras are fixed statically with the background.

143 Artificial features, such as fiducial markers, have strong (black-white) contrast and sharp
 144 features with straight edges for robust detection compared to natural features. However, severe
 145 structural motion during a shake table test can make fiducial marker detection in blurred images a
 146 really challenging task. Motion blur can deteriorate the marker detection performance using image
 147 processing algorithms (e.g., edge detection, blob detection). Deblurring methods can render clearer
 148 images for accurate feature detection. The underlying problem of the image deblurring part in this
 149 study is to restore clearer and sharper visual features on images for accurate detection and
 150 displacement computation.

151 (2) Motion Estimation for Severe Blur

152 In practice, severe motion blur on images can cause misdetections even when deblurring
 153 technique is implemented. For example, when the structure is subjected to vibration with a
 154 frequency close to the natural frequencies, structural vibration becomes much severer making the
 155 top floor shaking faster than the other floors. A post processing of data fitting is needed to
 156 complement the measurement in these cases. Although linear interpolation and/or spline
 157 interpolation can be used as basic data fitting considering the continuous movement of structure
 158 using nearby measurements, these interpolation methods neglect the system information and
 159 sometimes can yield wrong estimates at the misdetection instances.

160 Assume a measurement from a shake table test is denoted as \mathbf{y}_k at instance $t_k (k = 1, 2, \dots, T)$
 161 and the corresponding state is denoted as \mathbf{x}_k . For example, the measurement includes displacement
 162 measurement in x direction ($\mathbf{y}_k = d_x(t_k)$) for a 1D shake table test and the state include both
 163 displacement and velocity, such as $\mathbf{x}_k = [d_x(t_k), \dot{d}_x(t_k)]$. To model the motion, a state vector
 164 $\mathbf{x}_k \in \mathbb{R}^{n \times 1}$ denotes the system state and the linear dynamic system can be expressed as:

$$\mathbf{x}_k = \mathbf{A}_{k-1} \mathbf{x}_{k-1} + \mathbf{q}_{k-1} \quad (3)$$

165 where $\mathbf{x}_k \in \mathbb{R}^{n \times 1}$ is the state (as a vector), $\mathbf{q}_{k-1} \in \mathbb{R}^{n \times 1}$ is the process noise with Gaussian
 166 probability distribution $\mathbf{q}_{k-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{k-1})$, and $\mathbf{A}_{k-1} \in \mathbb{R}^{n \times n}$ denotes the dynamic
 167 model transition matrix.

168 The measurement equation is:

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{x}_k + \mathbf{r}_{k-1} \quad (4)$$

169 where $\mathbf{y}_k \in \mathbb{R}^{m \times 1}$ is the measurement, $\mathbf{r}_k \in \mathbb{R}^{m \times 1}$ is the measurement noise with Gaussian
 170 probability distribution $\mathbf{r}_{k-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_{k-1})$, and $\mathbf{H}_k \in \mathbb{R}^{m \times m}$ denotes the measurement
 171 model/matrix.

172 Successful observations are denoted as $(\mathbf{y}_k)_i \equiv (\hat{\mathbf{y}}_k)_i, (k, i) \in \mathcal{K}$:

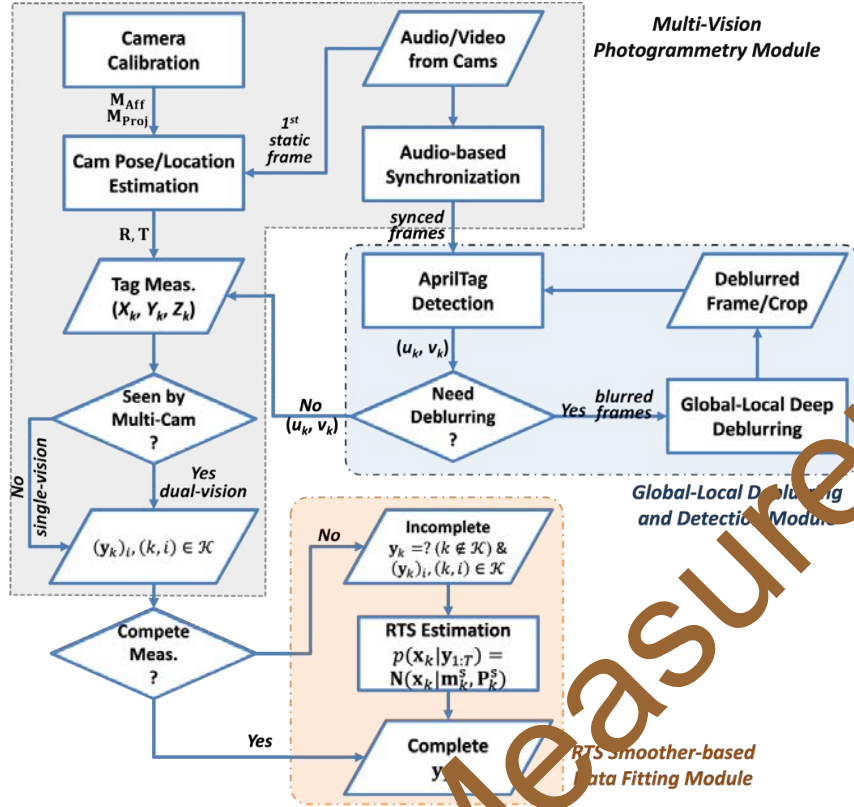
$$(\mathbf{y}_k)_i = (\mathbf{H}_k \mathbf{x}_k + \mathbf{r}_{k-1})_i \quad (k, i) \in \mathcal{K} \quad (4)$$

173 where $\mathcal{U} = \{1, 2, \dots, T\} \times \{1, 2, \dots, m\}$ denotes the universal set of scalar outputs corresponding to
 174 all possible observations, $\mathcal{K} \subseteq \mathcal{U}$ denotes the set corresponding to successful observations, \mathcal{M}
 175 denotes the set corresponding to failed/missed observations. $\mathcal{K} \cap \mathcal{M} = \emptyset$ and $\mathcal{K} \cup \mathcal{M} = \mathcal{U}$.

176 Severe structural vibrations sometimes can make motion blur so blemished that CV-based
 177 measurements cannot be obtained successfully even after using some image deblurring techniques.
 178 For discussion, these failed/missed observations are denoted as $(\mathbf{y}_k)_i \equiv (?_k)_i, (k, i) \in \mathcal{M}$. The
 179 underlying mathematical problem of the data fitting part in this study is to estimate failed/missed
 180 measurements $(?_k)_i, (k, i) \in \mathcal{M}$ based on the successful observations, $(\hat{\mathbf{y}}_k)_i, (k, i) \in \mathcal{K}$.

181 3. Method

182 Motion blurs involve shutter speed of camera and moving speed of recorded object. The
 183 failure in feature detection resulting from motion blur can cause misdetection of key features on
 184 vibrating structures, leading to empty observations at certain time instances. Hence, it is important
 185 to develop a deblurring and detection strategy that is suitable for displacement monitoring in shake
 186 table tests. Targeting toward the motion blur issue in shake table tests, a multi-vision approach was
 187 proposed including three modules (as shown in **Figure 2**): multi-vision photogrammetry module,
 188 global-local deblurring and detection module, and data fitting module for severe motion blurs.



189
 190 **Figure 2.** Flow chart of the proposed approach for displacement monitoring with multi-vision photogrammetry
 191 module, image deblurring and detection module, and data fitting module.

192

193 3.1. Multi-Vision Displacement Monitoring

194 In a shake table test, there are situations that not all tags are in the scope of view. For example,
 195 the top floor of a structure might be out of the camera due to excessive displacements, or one
 196 feature for tracking may be blocked by a structure component from one camera view. Sometimes
 197 during experiments, visual features on structures (e.g., artificial patterns, natural features) cannot
 198 be seen clearly or easily due to limited conditions (e.g., poor illumination, unsatisfactory camera
 199 pose, large movement of structure). To cope with these non-perfect situations in practice, a multi-
 200 vision strategy with both single-vision and dual-vision choices is needed to obtain full-field
 201 measurement in the post data analysis. In addition, vibrating structure may induce different levels
 202 of blurs viewed in multiple perspectives. Even if motion blur (**Figure 1a**) is too severe to be viewed
 203 clearly in one perspective, it doesn't necessarily mean that the blur will be at the same level with
 204 another perspective. The effectiveness of image deblurring of the same visual features may differ

205 due to different perspectives. Therefore, multi-vision scheme was chosen to provide information
206 redundancy for dynamic experiments.

207 ***(1) Fiducial Marker Detection and Video Synchronization***

208 Visual features in images used for displacement measurements could either be natural features
209 (e.g., structural corners) or artificial features (e.g., fiducial markers). Marker-free methods require
210 no speckle pattern or marker deployment, but they require more computation time to process
211 images to get features for matching. In contrast, artificial markers are commonly used to obtain
212 streamlined detection and tracking of points of interest (Spencer Jr et al., 2019). In order to process
213 videos in a fast manner, a fast feature detection and association (across camera) algorithm is
214 preferred for shake table tests. Therefore, fiducial markers (e.g., AprilTag (Olson, 2011)) with
215 sharp features was on top of the list for this study, as well as a speedy detection algorithm. AprilTag
216 detection algorithm (Wang & Olson, 2016) includes the first step of quad detection and the second
217 step of detailed pattern recognition. In the second step, a quad candidate generated from the first
218 step will be decoded to compare with the tag dictionary in the family to decide if the binary payload
219 matches with one specific tag pattern. This work aims to address the issue of motion-induced blur
220 in shake table tests and the proposed method can be integrated with different types of visual
221 features for CV-based monitoring. For demonstration purpose, AprilTag's were used as example
222 for feature detection.

223 In the first step of tag detection, quad detection may fail if line/quad features are blemished
224 by motion blur. In the second step, even if a marker is detected as a candidate quad in the first step,
225 the decoder will filter a marker out if its binary pattern is wrecked by motion blur, leading to no
226 match in the known tag family (Krogius et al., 2019; Liu et al., 2022). Therefore, there is a need
227 to restore images and recover the sharp features of the markers before achieving tag detection.

228 Followed by structural feature detection, video synchronization is of great importance for
229 multi-vision application, especially for shake table tests. One may argue to have all the cameras are
230 triggered at the same time in the beginning of the shake table tests to enforce video frames match.
231 However, the internal clock within each of the cameras will slightly drift during the recording
232 (especially for non-expensive cameras), making a mismatch of video frame across different
233 cameras. The mismatched frames will yield considerable error in multi-vision triangulation
234 computation. Therefore, a post video synchronization is needed before the image processing. In
235 this study, the ambient sound from shake table tests was recorded on the audio channels and the

236 audio recordings were processed and synchronized based on audio waveform matching using
 237 cross-correlation.

238 (2) Multi-Vision Triangulation

239 In a dual-vision setting, images from two different perspectives (**Figure 1b**) can serve as
 240 strong constraint in 3D scene reconstruction when the two viewing rays corresponding to the same
 241 scene point intersect. The 3D coordinates can then be determined using the direct linear transform
 242 (DLT) method (Abdel-Aziz et al., 2015) based on triangulation.

$$s\mathbf{X}_{p'} = \begin{bmatrix} A^{T^k w} \mathbf{X}_P \\ B^{T^k w} \mathbf{X}_P \\ C^{T^k w} \mathbf{X}_P \end{bmatrix} \quad (6)$$

243 where $A^{T^{(k)}}$, $B^{T^{(k)}}$ and $C^{T^{(k)}}$ represent the three rows of the transformation matrix $\mathbf{M}_{\text{trans}}^{(k)}$ for k -
 244 th camera. Transformation matrix $\mathbf{M}_{\text{trans}} = \mathbf{M}_{\text{aff}} \mathbf{M}_{\text{proj}} [\mathbf{R} | \mathbf{T}]$.

245 In this study, pinhole model (Eq. 1) was used for each of the two points in the two images
 246 (e.g., p'_1, p'_2 in **Figure 1b**), respectively.

$$\begin{cases} s x_{p'_i} = \frac{s u_{p'_i}}{s w_{p'_i}} = \frac{A^{T^k w} X_P}{C^{T^k w} X_P} \\ s y_{p'_i} = \frac{s v_{p'_i}}{s w_{p'_i}} = \frac{B^{T^k w} X_P}{C^{T^k w} X_P} \end{cases} \rightarrow \begin{cases} (s x_{p'_i} C^{T^k} - A^{T^{(k)}}) w X_P = 0 \\ (s y_{p'_i} C^{T^k} - B^{T^{(k)}}) w X_P = 0 \end{cases} \quad (7)$$

247 Combining the equations developed from the two points (Eq. 7), linear algebra equations are
 248 derived to yield a unique solution. The four observations ($s u_{p'_i}$ and $s v_{p'_i}$ from each point) make
 249 it a determinate problem to solve.

$$\begin{bmatrix} s x_{p'_1} C^{T(1)} - A^{T(1)} \\ s y_{p'_1} C^{T(1)} - B^{T(1)} \\ s x_{p'_2} C^{T(2)} - A^{T(2)} \\ s y_{p'_2} C^{T(2)} - B^{T(2)} \end{bmatrix}_{4 \times 4} \begin{bmatrix} w U_P \\ w V_P \\ w W_P \\ w T_P \end{bmatrix}_{4 \times 1} = \mathbf{0} \quad (8)$$

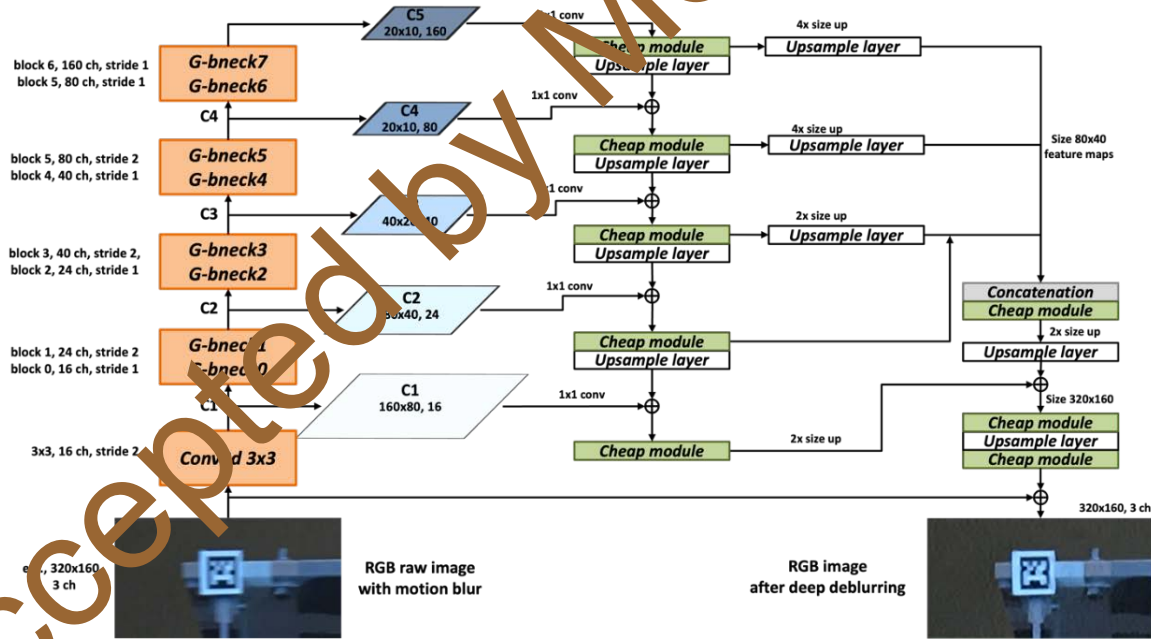
250 where $w T_P$ is treated as a scale factor and the homogeneous coordinates of point P could be
 251 represented as $[w X_P, w Y_P, w Z_P, 1]^T$.

252 3.2. Global-Local Detection and Deblurring (GLDD)

253 Restoring images from local blurring of moving objectives is an open problem and image-
 254 based deep deblurring methods are yet to be applied in CV-based monitoring of shake table tests.
 255 One focus of this study lies on the design of global-local detection and deblurring (GLDD) module
 256 using deep deblurring to augment the displacement measurement in shake table tests.

257 **(1) Deep Deblurring Model**

258 In the study, a deblurring model was adopted that is a generative adversarial network (GAN)-
 259 based deep deblurring model. GANs were investigated for image restoration (Pramakrishnan,
 260 2017) by refereeing to the idea of image translation and the recent development includes
 261 DeblurGAN (Kupyn et al., 2018), DeblurGAN v2 (Kupyn, 2019), and Ghost-DeblurGAN (Liu et
 262 al., 2022). Due to the high speed of processing, models with light-weight convolution neural
 263 network (CNN)-based feature extractors, such as GhostDeblurGAN, are preferred in the study for
 264 efficient image deblurring compared with heavyweight models. Hence, this study adopted a
 265 lightweight deblurring model, DeblurGAN-v2, as the image restoration component for the
 266 proposed deblurring module.



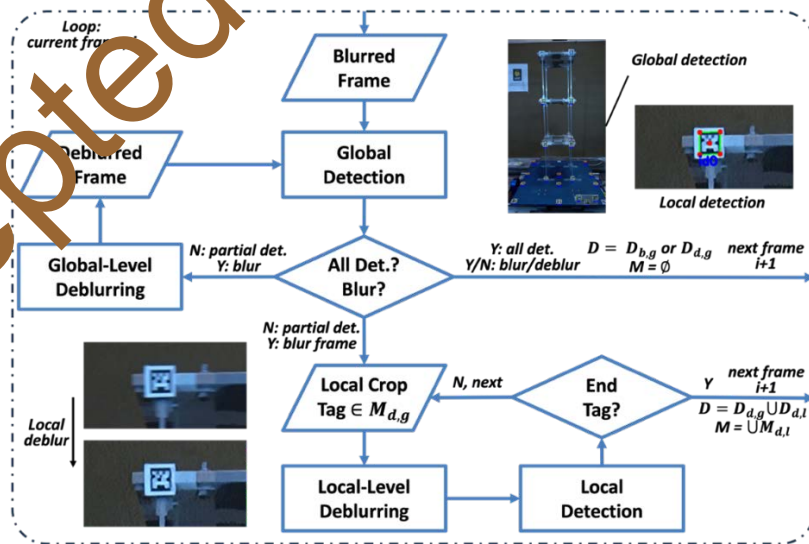
267 **Figure 3.** Schematic view of the Ghost-DeblurGAN generator architecture with an example of image
 268 processing on fiducial marker attached to a structure.
 269

270 The generator part of a GAN-based deep deblurring model includes a CNN backbone as
 271 feature extractor and a feature pyramid network for a rich set of global feature maps at different
 272 spatial scales (**Figure 3**). Compared to DeblurGAN-v2, the CNN backbone of GhostDeblurGAN

273 uses GhostNet (Han, 2020) as shown in **Figure 3** instead of MobileNet. Cheap module (Sandler,
 274 2018) adopted in the GhostNet CNN backbone includes the pointwise and depthwise separable
 275 convolution layers in sequence to obtain the intrinsic feature maps. These intrinsic feature maps
 276 are comparable to computational expensive 2D convolution layers, but the yielding has a much
 277 faster speed. Compared to the 2D convolution layers in DeblurGAN-v2, the cheap modules in
 278 Ghost-DeblurGAN can help reduce 53.21% of the floating-point operations per second (FLOPS),
 279 during the forward calculation. Hence, Ghost-DeblurGAN model was adopted as the deblurring
 280 model of the proposed framework. The deblurring model was trained using a large scale YorkTag
 281 dataset (Liu et al., 2022) with paired blurry-sharp tag images of 2074 pairs (1577 in training set
 282 and 497 in test set) that were collected in both indoor and outdoor environments. The deblurred
 283 images (9761 tags within dataset) processed using the trained deblurring model showed an
 284 improved detection rate of 59.1%, compared to detection rate of only 32.0% when using raw
 285 blurred images.

286 **(2) Global-Local Detection and Deblurring (GLDD) Module**

287 A global detection and deblurring (GDD) process on the whole image may be sufficient to
 288 restore images from small motion blurs. However, when the motion blur becomes severer, local
 289 level of deblurring process is needed on the key locations to retrieve sharp visual clue for tag
 290 detection. Therefore, global-level and local-level image deblurring were automated to augment
 291 tag detection based on the different extent of motion blur.



292

293 **Figure 4.** Flow chart of the proposed GLDD module. (note: both artificial and natural features can be
 294 used and AprilTag is adopted as an example)

295 The proposed GLDD module (**Figure 4**) includes global-level tag detection/deblurring on
 296 whole video frames and local-level detection/deblurring on cropped images near the key (image)
 297 locations. The assumption for such design is that the attention of a general deep deblurring model
 298 is distracted on nonimportant area when dealing with a relatively large moving foreground (e.g.,
 299 the vibrating structure) instead of focusing on the key structural features. If an image crop contains
 300 single features/tags with relative larger foreground, the attention of a CNN-based feature extractor
 301 will be forced to put on the tags over the background and the tags will be less difficult to restore.
 302 For a shake table test, rigid movement occurs at the sliding base. If the structure moves at the same
 303 frequency as the sliding base when subjected to a forced vibration, the moving distance on the
 304 higher levels of the structure over the same time would be larger compared to the sliding base.
 305 Hence, it is reasonable to assume that the blur severity at structure top is larger than the base.

306 **3.3. Data Fitting for Severe Motion Blur**

307 The problem of data fitting was reshaped with a perspective of Kalman filtering (KF) (Welch
 308 & Bishop, 1995) and Rauch-Tung-Striebel (RTS) smoothing (Särkkä, 2008) thinking. In this
 309 study, the dynamic system was described by a partially observed Markov process in the Bayesian
 310 sense by computing the conditional distributions (e.g., $p(\mathbf{x}_k|\mathbf{x}_{k-1})$) using either filtering or
 311 smothering methods (Särkkä & Svensson, 2023).

313 **(1) Kalman Filter-based Estimation**

314 KF can estimate current state of a dynamical system (\mathbf{x}_k) given previous and current
 315 observations ($\mathbf{y}_j, j = 1, 2, \dots, k$). KF probabilistic state model consists of the conditional
 316 probability distributions of the state and the measurement which are Gaussian distributions:

$$\mathbf{x}_k \sim p(\mathbf{x}_k|\mathbf{x}_{k-1}) = \mathbf{N}(\mathbf{x}_k|\mathbf{A}_{k-1}\mathbf{x}_{k-1}, \mathbf{Q}_{k-1}) \quad (9)$$

$$\mathbf{y}_k \sim p(\mathbf{y}_k|\mathbf{x}_k) = \mathbf{N}(\mathbf{y}_k|\mathbf{H}_k\mathbf{x}_k, \mathbf{R}_k) \quad (10)$$

317 In order to compute the filtering results, the parameters or states are computed in two steps
 318 recursively: prediction step and correction step. The prediction step in the recursive computation
 319 includes the mean prediction and covariance prediction:

$$\mathbf{m}_k^- = \mathbf{A}_{k-1} \mathbf{m}_{k-1} \quad (11)$$

$$\mathbf{P}_k^- = \mathbf{A}_{k-1} \mathbf{P}_{k-1}^- \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1} \quad (12)$$

320 The difference between the predicted measurement and the sensor reading is denoted as
 321 innovation \mathbf{v}_k :

$$\mathbf{v}_k = \mathbf{y}_k - \mathbf{H}_k \mathbf{m}_k^- \quad (13)$$

322 The correction step in the recursive computation includes the correction of mean and
 323 covariance of the current state:

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \quad (14)$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T \quad (15)$$

324 where $\mathbf{m}_k^-, \mathbf{m}_k \in \mathbb{R}^{n \times 1}$ are the mean value of the state \mathbf{x}_k and $\mathbf{P}_k^-, \mathbf{P}_k \in \mathbb{R}^{n \times n}$ are the covariance
 325 matrix of the measurement \mathbf{x}_k during the prediction and correction steps, respectively. Kalman
 326 gain for the correction is $\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T \mathbf{S}_k^{-1}$. Covariance matrix for the innovation is $\mathbf{S}_k =$
 327 $\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k$.

328 From the previous state \mathbf{x}_{k-1} and current measurement \mathbf{y}_k , a prediction and correction can be
 329 performed to estimate current state of \mathbf{m}_k and \mathbf{P}_k sequentially. The next prediction of the
 330 measurement can be obtained as $\mathbf{H}_{k+1} \mathbf{m}_{k+1}^-$ if there is a misdetection of \mathbf{y}_{k+1} , $k+1 \notin \mathcal{K}$.
 331 However, there is no correction step at $k+1$ due to the lack of measurement which can lead to
 332 growing covariance matrix \mathbf{P}_{k+1} .

333 (2) Smoother-based Estimation

334 RTS smoother is the smoothing method of estimating the current state given the whole
 335 measurements instead of just using the current measurement and the previous state. Because on-
 336 time measurement is not required in a shake table test, a short-time delay (few seconds or minutes)
 337 is allowed and can be used for post processing. Therefore, the study considered using the available
 338 measurement not that were not just prior to the current steps ($1 \leq i \leq k$ & $i \in \mathcal{K}$) but also after
 339 the current steps ($k+1 \leq i \leq T$ & $i \in \mathcal{K}$). RTS smoother, which is similar to but not same as the
 340 backward algorithm of KF, is used to estimate the missed measurement $\mathbf{y}_k = ?$ ($k \notin \mathcal{K}$) using more
 341 measurement in future ($\mathbf{y}_i, k+1 \leq i \leq T$ & $i \in \mathcal{K}$) in addition to KF. In the study, the KF and
 342 RTS smoother results would be combined with expectation-maximization to estimate the dynamic
 343 state and missing measurement in the shake table tests. Unlike normal RTS smoother that uses all

344 the measurements for all time steps, this study would impose a constraint on measurement because
 345 only partial measurements ($\mathbf{y}_i, 1 \leq i \leq T \ \& \ i \in \mathcal{K}$) could be provided for smoothing.

346 The close form smoothing solution for a RTS smoother is:

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}) = \mathbf{N}(\mathbf{x}_k | \mathbf{m}_k^s, \mathbf{P}_k^s) \quad (16)$$

347 where $\mathbf{m}_k^s, \mathbf{P}_k^s$ are the estimated mean and covariance of the current state \mathbf{x}_k based on the whole
 348 measurements $\mathbf{y}_{1:T}$.

349 RTS smoother allows one to refine estimates of current states using the information provided
 350 by later observations. The equations for the backward recursion for RTS smoother include the
 351 prediction step:

$$\mathbf{m}_{k+1}^- = \mathbf{A}_k \mathbf{m}_k \quad (17)$$

$$\mathbf{P}_{k+1}^- = \mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T + \mathbf{Q}_k \quad (18)$$

352 where $\mathbf{m}_k \in \mathbb{R}^{n \times 1}$ and $\mathbf{P}_k \in \mathbb{R}^{m \times m}$ are the mean value and the covariance matrix of the state \mathbf{x}_k
 353 computed by the KF.

354 The correction step in the recursive computation includes:

$$\mathbf{m}_k^s = \mathbf{m}_k + \mathbf{G}_k (\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-) \quad (19)$$

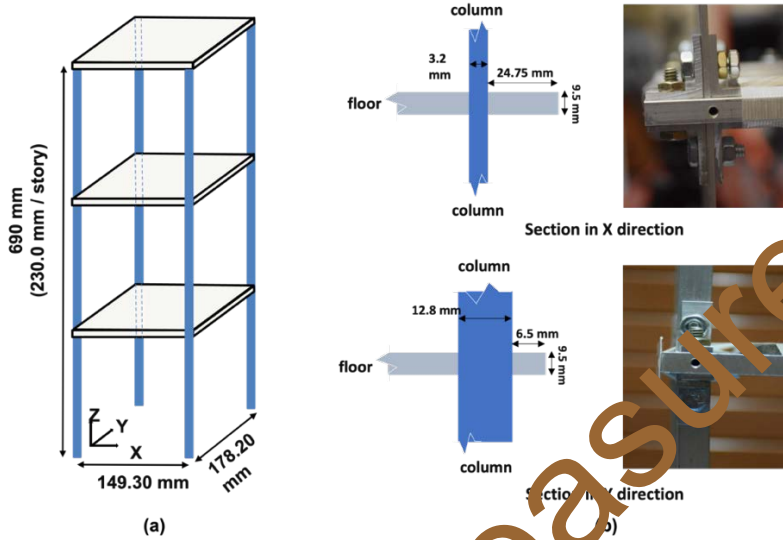
$$\mathbf{P}_k^s = \mathbf{P}_k + \mathbf{G}_k (\mathbf{P}_{k+1}^- - \mathbf{P}_{k+1}^-) \mathbf{G}_k \quad (20)$$

355 where $\mathbf{m}_k^s \in \mathbb{R}^{n \times 1}$ and $\mathbf{P}_k^s \in \mathbb{R}^{m \times m}$ are the mean and the covariance matrix of the measurement
 356 \mathbf{x}_k computed by the RTS smoother. $\mathbf{G}_k = \mathbf{P}_k \mathbf{A}_k^T (\mathbf{P}_{k+1}^-)^{-1}$ is the smoother gain for the correction.

357 4. Experiment: Shake Table Test

358 To evaluate the proposed multi-vision approach and algorithm for displacement monitoring, a
 359 shake table test was carried out on a three-story aluminum frame (**Figure 5**). Chirp excitation was
 360 used as the input ground excitation to induce different levels of structural responses. The three-
 361 story aluminum frame (**Figure 5a**) were fabricated with the same story heights of 230 mm. The
 362 width and length of the floors in X and Y directions are 202 mm and 204 mm, respectively. The
 363 detailed views of the column to floor connection are shown in **Figure 5b** with X direction as the
 364 weak direction and Y direction as the strong direction. The center-to-center distances between the
 365 two adjacent columns are 149.30 mm in X direction and 178.2 mm in Y direction, respectively. A
 366 steel plate with the same mass, 0.66 kg, was affixed to the center of each floor. A shake table
 367 (Quanser Shake Table II) was utilized to provide lateral excitation with a payload area size of 460
 368 mm \times 460 mm. The maximum stroke limit of the actuator is ± 76.2 mm and the frequency range

369 of input motion is 0.5-10 Hz. The proposed approach and LVDT were used to measure dynamic
 370 displacements. Numerical simulation, experimental setup, and feature detection (without GLDD
 371 implementation) are presented as follows.



372
 373 **Figure 5.** (a) 3D schematic view of the aluminum frame and (b) detailed views of the column-floor
 374 connection in X and Y directions.

375 **4.1. Finite Element Simulation**

376 A finite element (FE) model of the same aluminum frame was developed and analyzed using
 377 OpenSees (Mazzoni et al., 2006) to understand the structural behavior of the physical model. The
 378 same geometry (**Figure 5**) was used to design the FE model and the columns and floors are
 379 modeled by assigning fiber sections to dispBeamColumn and ShellMITC4 elements in OpenSees,
 380 respectively. The mechanical properties of the aluminum material for the simulation were: yield
 381 strength = $2.5e8$ N/m², modulus of elasticity = $6.9e10$ N/m², Poisson's ratio = 0.33, density = 2700
 382 kg/m³. Following the experimental model, lumped masses of 0.66 kg were assigned to all the three
 383 stories. The first three modal frequencies in the X direction of the FE model were 5.68 Hz, 16.06
 384 Hz, and 23.18 Hz (see **Table I**) based on the modal analysis. Ground excitations using an upchirp
 385 excitation (designed as 0.5-4.5 Hz) and the free vibration after the excitation are simulated on the
 386 FE model. By knowing the modal parameters (e.g., modal frequencies) from FE analysis, the
 387 ground excitation can be well designed for the real experiment to cover different frequency
 388 spectrums while maintaining the safety during the laboratory test.

389

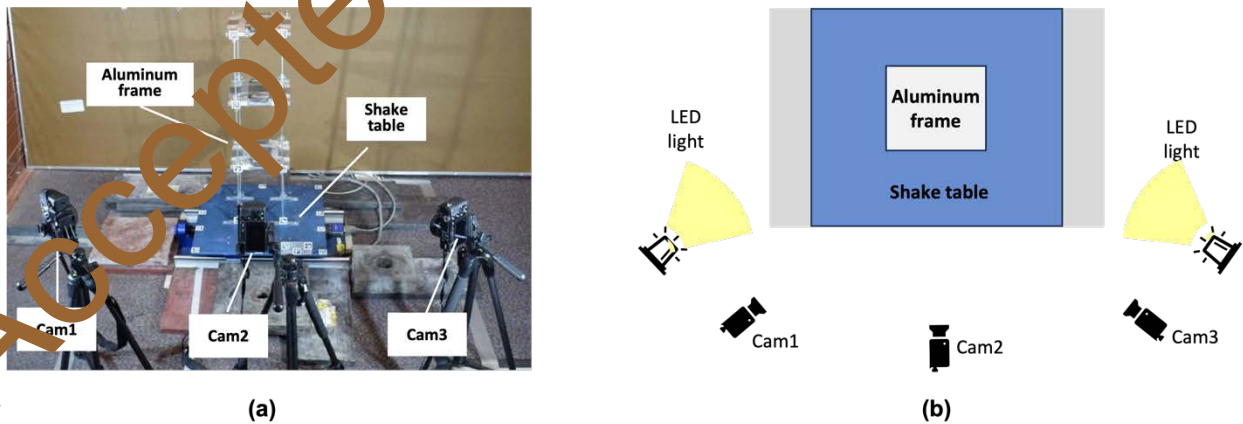
390 **Table 1.** Natural frequency (first three modes in X direction) and MAC value comparison between FEM
 391 and output-only system identification results using (virtual) free vibrations and experimental measurements.

Meas.	Sys. ID Method	Mode1		Mode2		Mode3	
		Freq. (Hz)	MAC	Freq. (Hz)	MAC	Freq. (Hz)	MAC
Virtual (FEA)	Modal Analysis	5.68	-	16.06	-	23.48	-
	FDD	5.62	1.00	15.87	1.00	23.19	1.00
	SOBI	5.62	1.00	15.87	0.96	23.07	0.93
	SSI	5.60	1.00	15.85	1.00	23.09	1.00
Experiment (CV)	FDD	5.15	-	14.95	-	24.55	-
	SOBI	5.15	-	14.99	-	24.59	-
	SSI	5.14	-	15.02	-	24.55	-

392

393 **4.2. Experimental Setup**

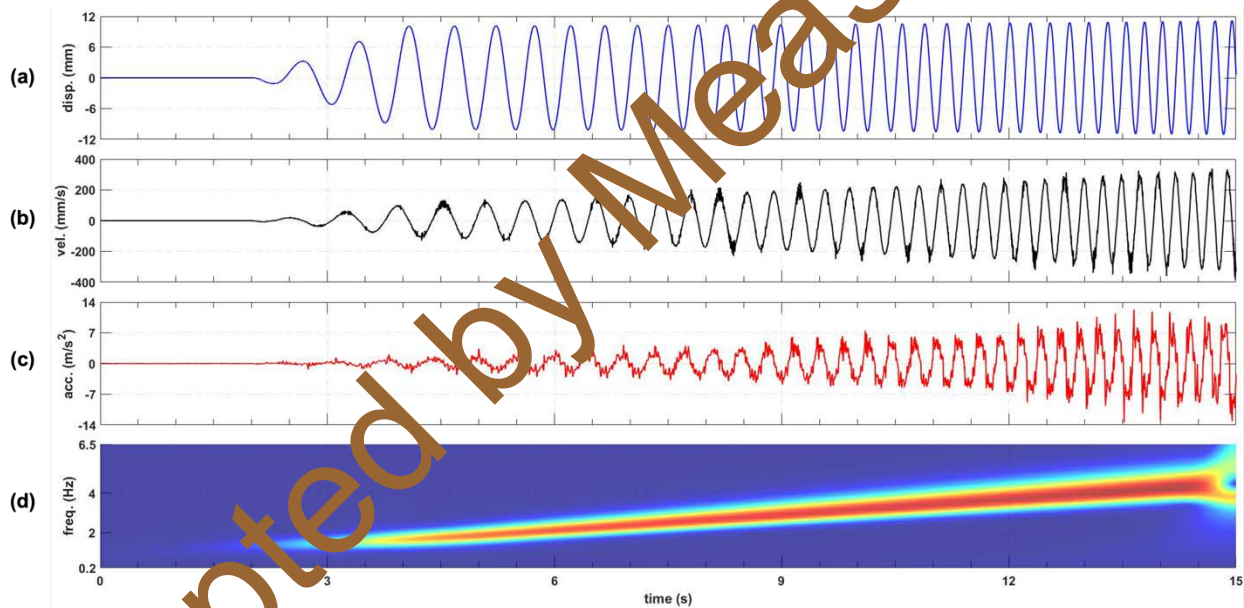
394 Twenty-three AprilTags of “25h9 tag family” with unique IDs were attached to the aluminum
 395 frame. Eight tags were attached to the key end locations on the front surface of frame’s floors to
 396 record displacement time histories of the structure during the experiment. The remaining seventeen
 397 tags were attached onto the surface of the table base to perform camera location/pose estimation
 398 and to record the displacement time history of the base. Moreover, the displacement and
 399 acceleration time histories of the base were recorded by the LVDT and accelerometer integrated
 400 with the shake table.



401

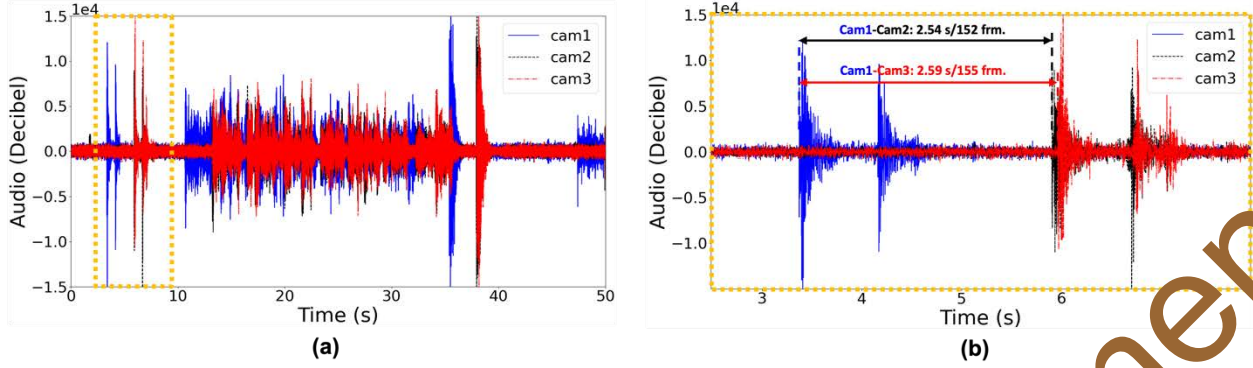
402 **Figure 6.** (a) Photo and (b) schematics of the experimental setup of the shake table test with the aluminum
 403 frame and cameras.

404 The aluminum frame (**Figure 6a**) was subjected to the same upchirp (0.5-4.5 Hz) excitation
 405 (**Figure 7**) as in the FE simulation (**Figure 5**) with the maximum base displacement of ± 10 mm
 406 and the duration of 15 seconds. Three portable cameras (**Figure 6**) were used to monitor the shake
 407 table tests: two Sony $\alpha 6400$ cameras (denoted as Cam1 and Cam2) and one Sony $\alpha 6000$ camera
 408 (denoted as Cam3). Multiple DC light-emitting diode (LED) sources were deployed further away
 409 from the shake table for balanced illumination condition. Please note that there is a trade-off
 410 between different (camera) angles of view: wide-angle allows more context in view while
 411 sacrificing the density of pixels over foreground (structures); narrow angle allows denser pixels
 412 over structures while covers less area. Based on the size of the three-story frame structure, a
 413 medium field view ($f=16$ mm) was adopted with field of view angles of 72.59° in vertical and
 414 52.27° in horizontal directions. Camera parameters were set as the same for the three cameras to
 415 record high-quality (1920×1080 pixel²) videos of the structure in vibration: focal length = 16 mm,
 416 ISO = 2000 (the sensitivity to light), frame rate = 59.94 fps, shutter speed = $1/165$ s (6.1 ms), and
 417 camera aperture = f6.3.



418 **Figure 7.** Time histories of (a) the base displacement of the chirp ground excitation (measured by LVDT),
 419 (b) velocity using 1st order differentiation, (c) acceleration using 2nd order differentiation, and (d) the
 420 wavelet transform of the base displacement.
 421

422 Ambient soundtracks of the test were used for video matching between the three cameras
 423 using cross correlation method to compute the differences in time (**Figure 8**). Cam1-audio was
 424 used as a reference and the time shifts were +2.54 s (+152 frames) and +2.59 s (+155 frames) for
 425 the Cam2-audio (-video) and Cam3-audio (-video) channels, respectively. After the
 426 synchronization of frames from multiple perspectives, multi-vision triangulation is performed.

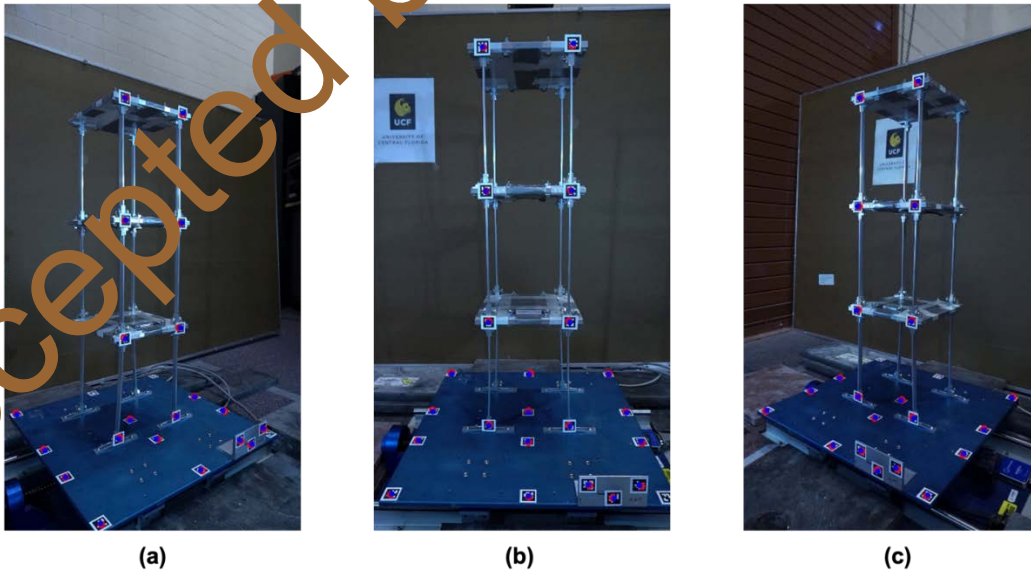


427

428 **Figure 8.** (a) Synchronization processing using the audio data from the three cameras and cross-
 429 correlation, and (b) the detailed view within the time window of (2.5-8 s).

430 4.3. Tag Detection without GLDD

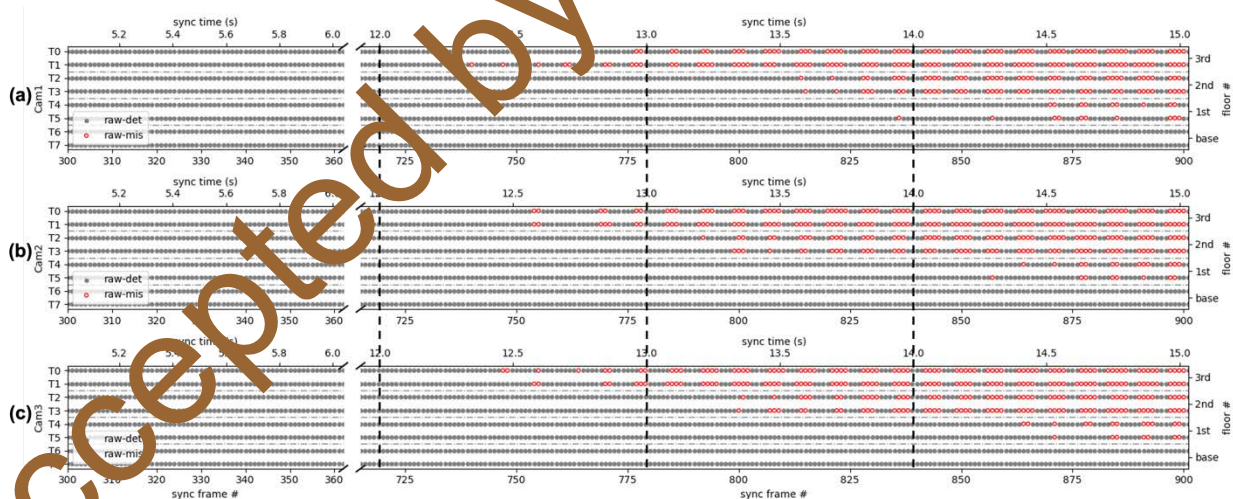
431 Evaluation of camera calibrations of the three cameras (**Figure 9**) were conducted with the
 432 blue dots representing the 3D location in WCS projected onto camera images (using recognized
 433 camera parameters) and red circles representing the detected location in SCS. The average error in
 434 2D SCS between the detection locations $[\hat{u}_i, \hat{v}_i]^T$ using AprilTags and the projected locations in
 435 2D SCS were 2.99 pixel for Cam1, 3.53 pixel for Cam2, and 2.97 pixel for Cam3, respectively.
 436 The average errors in 3D SCS between the computed location $[\hat{X}_i, \hat{Y}_i, \hat{Z}_i]^T$ and the measurements
 437 by rulers in 3D WCS were 2.03 mm for Cam1, 1.58 mm for Cam2, and 1.72 mm for Cam3,
 438 respectively.



439

440 **Figure 9.** Evaluation of camera calibration and pose estimation in the small-scale shake test. (note: red
 441 circles are projection using camera parameters and the ground truth locations in WCS and blue dots are
 442 detected points in SCS)

443 The detection performance is shown in **Figure 10** for the three cameras. A successful detection
 444 event by raw tag detection technique is denoted as a gray solid circle and a failed detection event
 445 is denoted as a red circle for each of the three perspectives at each synchronized frame/time across
 446 the three cameras. The detection performances are shown for different floors from top to bottom:
 447 T0/T1 on the 3rd floor, T2/T3 on the 2nd floor, T4/T5 on the 1st floor, and T6/T7 on the base. When
 448 the motion blur was little (0-12 s, 0-720 frames), the detection performance for the raw tag
 449 detection is satisfactory with all the tags on the frame successfully detected and the success rate
 450 was 100%. As the vibration becomes large enough to cause mild motion blur (12.3-13.2 s), the
 451 tags on the 3rd floor (T0 and T1) were difficult to identify. As the excitation frequency increased
 452 (13.2-14.2 s), tags on the 3rd floor more frequently failed to be detected and tags on the 2nd floor
 453 experience misdetections. During the last one second, when the excitation frequency was close to
 454 the 1st natural frequency of the structure, even tags on the 1st floor were difficult to be detected just
 455 using the raw tag detection technique. T6 and T7 on the base floor showed a 100% detection rate
 456 for all the cameras.



457 **Figure 10.** Tag detection evaluation based on raw image frames from (a) Camera 1, (b) Camera 2, and (c)
 458 Camera 3 in the shake table test.

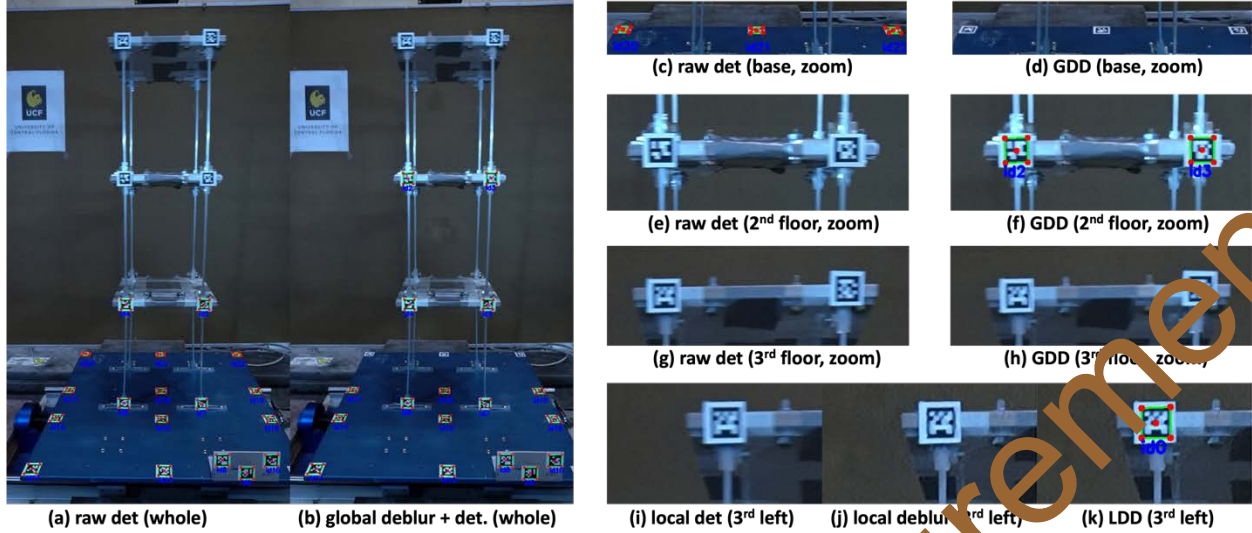
460 The vertical lines in **Figure 10** delineate the time segments which are associated with different
 461 degrees of motion blur. As shown in **Table 2**, the missed rates for the mild level of motion blur
 462 (12-13 s) were 12/480, 13/480, and 13/480 for Cam1-3, respectively. The missed rates for the

463 medium level of motion blur (13-14 s) increased to 70/480, 91/480, and 94/480 for Cam1-3,
464 respectively. The missed rates for the severe level of motion blur (14-15 s) were 151/480 (Cam1),
465 162/480 (Cam2), and 159/480 (Cam3). Misdetection events occur when the velocities of the tags
466 were higher on upper floors. In order to obtain more dynamic measurements, GLDD is needed to
467 enhance the detection rate.

468 5. Result and Discussion

469 5.1. Augmented Detection with Multi-Vision and GLDD

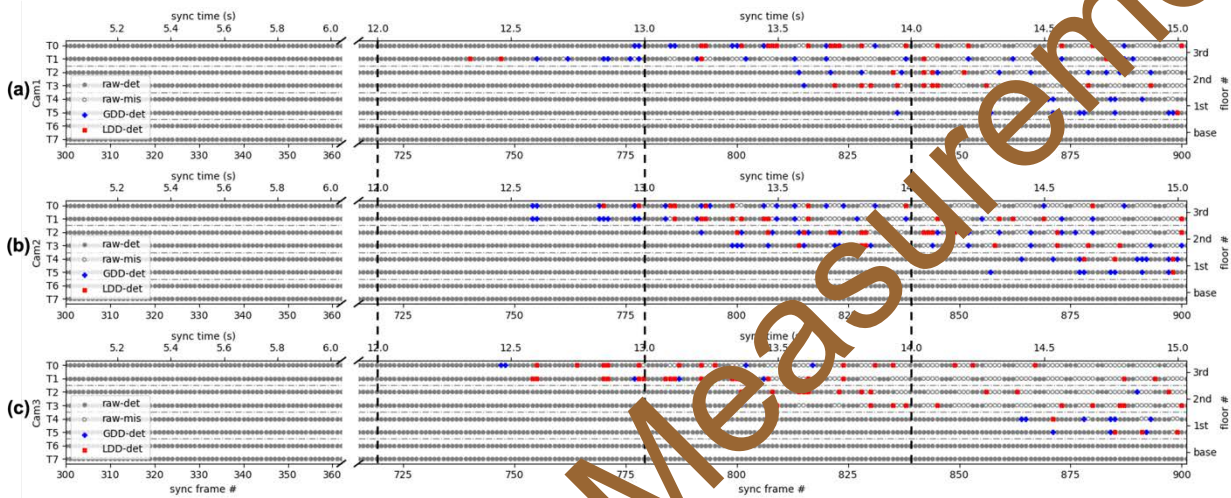
470 As an example, the image-based GLDD process on a video frame is shown in **Figure 11**.
471 When there was motion blur on images, the detection algorithm using the raw blurred image
472 (**Figure 11a**) could not identify all the tags on the structure, in contrast to a static frame without
473 motion blur (**Figure 9b**). However, the detection rate can be improved by using the GLDD module
474 to an extent. For example, the GDD algorithm could improve the detection on the front surface of
475 the frame (**Figure 11b**), especially on the 2nd floor (**Figure 11f**) that was at the center of the whole
476 image. However, while the GAN-based deblurring on the whole image could only improve
477 detection in a certain focused area (e.g., center), the global image restoration was not enough to
478 make the tags on the 3rd floor detectable (**Figure 11h**). In addition, the objects that were off
479 centered (e.g., base) did not necessarily become sharper using the GAN-based deep deblurring
480 method, which could even impose adverse effects making previously detected tags less detectable
481 (**Figure 11d**) after the image restoration. However, the local deblurring on a local crop image
482 (**Figure 11j**) could shift more attention on the important features (e.g., tags, structural features)
483 making the bit features sharper for successful detection (**Figure 11j**). One may argue that this
484 improvement might be due to the cropped image size. Experiment result (**Figure 11i**), though,
485 proved that using an image crop even around a tag would not necessarily improve success rate.



486
 487 **Figure 11.** GLDD process steps on one video frame: (a) detection result on raw image, (b) detection result
 488 on globally deblurred image, (c) raw detection result and (d) GDD processed result on the table base, (e)
 489 raw detection results and (d) GDD processed result on the 2nd floor, (g) raw detection result and (h) GDD
 490 processed result on the 3rd floor, (i) detection result on the local image near Tag0, (j) its locally deblurred
 491 image, and (k) LDD processed result.

492 The limit of the proposed framework was evaluated using the upchirp excitation whose actual
 493 frequency window is about 0.5-6.2 Hz covering the 1st natural frequency of the structure (5.15 Hz
 494 in **Table 1**) near the end of the experiment. Because the transient displacement input on the table
 495 base included frequency component at the 1st natural frequency near 15 s (**Figure 7**), the short-
 496 time resonance caused the frame to shake violently and resulted in severe motion blurs in videos.
 497 As shown in **Figure 12**, the detection performance for the GLDD process for each of the camera
 498 is presented with blue-plus symbols denoting successful detections using the GDD and red-cross
 499 symbols denoting successful detection using the local detection and deblurring (LDD). The
 500 detailed performances of GDD and LDD were compared using an ablation study of the different
 501 augmentation strategies (**Table 2**). When the time was 12-13 s in the shake table test, there were
 502 12 misses, 13 misses, and 13 misses for Cam1-3, respectively. With the GDD, the miss counts
 503 went down to 4, 2, and 10 for the three cameras, respectively. With the additional LDD, the miss
 504 counts went further down to 2, 0, and 0 for the three cameras, respectively, making the total
 505 restoration rates of 10/12 (83.3%) for Cam1, 13/13 (100.0%) for Cam2, and 13/13 (100.0%) for
 506 Cam3. The average restoration rate for the three cameras were 35/38 (92.1%) for mild-level motion
 507 blur. When the excitation frequency increases from around 3.9 Hz to 4.2 Hz during 13-14 s (**Figure**
 508 **7d**), the restoration rates for the GLDD process were 36/70 (51.4%) for Cam1, 54/91 (59.3%) for

509 Cam2, and 39/94 (41.5%) for Cam3. The average restoration rate for the three cameras was
 510 129/255 (50.6%) for the medium-level motion blur. During the last one second (14-15 s) when
 511 frequency span of the excitation (3.9-6.2 Hz) overlapped with the 1st natural frequency of the
 512 structure (5.15 Hz), the restoration rates by GLDD decreased to 44/151 (29.1%) for Cam1, 49/162
 513 (30.2%) for Cam2, and 26/159 (16.4%) for Cam3. The average restoration rate for the three
 514 cameras was 119/472 (25.2%) for severe-level motion blur.



515
 516 **Figure 12.** Tag detection evaluation with GLDD processing from (a) Camera 1, (b) Camera 2, and (c)
 517 Camera 3 in the shake table test.

518 After analyzing all the video frames during the whole shake table test (0-15.3 s, 0.5-4.5 Hz),
 519 it was found that the restoration rates of the GLDD were 92/243 (37.9%) for Cam1, 126/281
 520 (44.8%) for Cam2, 82/284 (28.9%) for Cam3. The different performances among cameras were
 521 due to the relative location and pose of cameras with respect to the shaking aluminum frame
 522 showing the effect of the camera placement on achieving high quality CV-based results. The
 523 GLDD process itself restored 94/207 (45.4%) of previous misdetections using raw images. The
 524 multi-vision strategy did take effect in the detection augmentation as well. Take Cam3 for
 525 example, the missed counts were brought from 13 down to 6 during 12-13 s, from 94 down to 60
 526 during 13-14 s, and from 159 to 133 during 14-15 s. For the whole test (0-15.3 s), the total miss
 527 count for Cam3 was brought from 284 down to 207 with a 26.0% drop. With the implementation
 528 of both strategies (multi-vision and GLDD), the total miss count for the experiment is brought
 529 down to only 113 (**Table 2**) with 75.0% measurements retrieved (from the previous misdetections)
 530 compared to just using raw images from one single camera (e.g., Cam3).

531

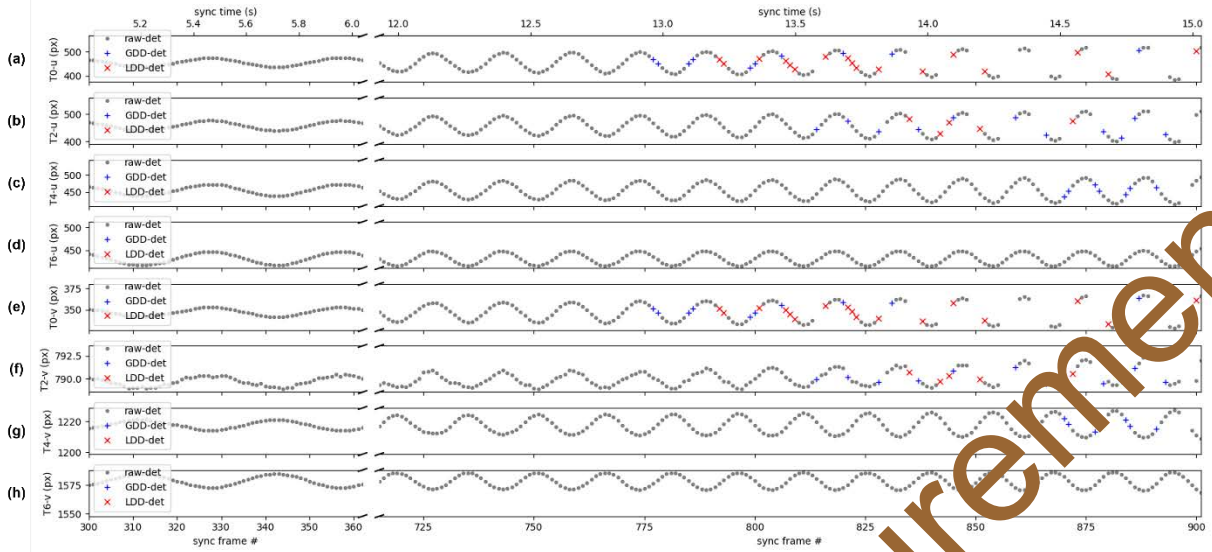
532 **Table 2.** Hit and miss counts using different tag detection methods for each camera and hybrid setting
 533 during different time windows of the shake table test.

detection method		t: 12-13s (frm: 719-778)				t: 13-14s (frm: 779-838)				t: 14-15s (frm: 739-899)				t: 0-15.3s (frm: 0-900)			
		cam1	cam2	cam3	multi	cam1	cam2	cam3	multi	cam1	cam2	cam3	multi	cam1	cam2	cam3	multi
raw	det	468	467	467	474	410	389	386	420	329	318	321	347	7133	7099	7099	7169
	mis	12	13	13	6	70	91	94	60	151	162	159	133	243	281	284	207
GDD	det	476	478	470	480	428	418	395	436	357	349	331	367	7187	7171	7115	7212
	mis	4	2	10	0	52	62	85	44	123	131	149	111	189	205	261	164
GLDD	det	478	480	480	480	446	443	425	461	373	367	347	390	7225	7221	7174	7263
	mis	2	0	0	0	34	37	55	19	107	113	133	90	151	155	202	113

534

535 **5.2. CV-based Displacement Measurement**

536 The pixel coordinates of the tag centers were localized on the collected video frames. For
 537 example, the image-based detection results (from Cam1) for the four floor levels (T0-3rd floor, T2-
 538 2nd floor, T4-1st floor, T6-base) are shown in **Figure 13**. Results from raw-image detection are
 539 denoted as gray dots, additional results from GDD process are denoted as blue “+” symbol, and
 540 the additional results from LDD process are denoted as red “×” symbol. From bottom to top of the
 541 structure, the increasing motion blur made it more and more difficult to detect using raw images.
 542 The GDD retrieved almost all the misdetections (denoted as “miss” in the study) on the 1st floor
 543 (**Figure 13c** and **g**) from 14.5-15.0 s (870-900 frames). On the 2nd floor and 3rd floor, the GLDD
 544 performed well from 12.5-13.8 s (775-825 frames) by restoring all the misdetections from raw
 545 images. LDD were more robust dealing with challenge events from 13.8-15.1 s (825–900 frames)
 546 when the motion blur was at a severe level.

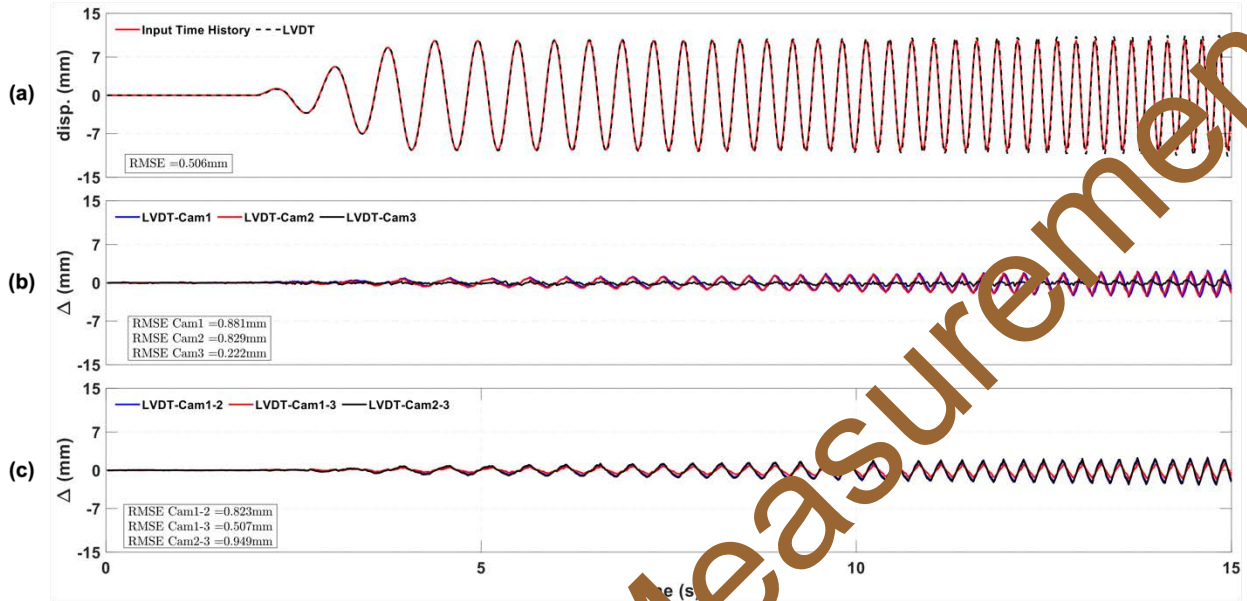


547
 548 **Figure 13.** Localization results of (a-c) u and (e-h) v sensor coordinate system for Tag0 on 3rd floor, Tag2
 549 on 2nd floor, Tag4 on 1st floor, and Tag6 on base using GLDD module

550 The time histories of displacement were compared (Figure 14) among the designed
 551 displacement input, LVDT measurement, single-vision measurements, and dual-vision
 552 measurements on the table base. There was a small difference (0.506 mm) between the designed
 553 (displacement) input and the LVDT measurement in the shake table test. This study used LVDT
 554 measurement as the baseline for comparison among vision-based measurements using root-mean-
 555 square error (RMSE). RSMEs between single vision methods and LVDT were 0.881 mm (Cam1),
 556 0.829 mm (Cam2), and 0.222 mm (Cam3), respectively. RSMEs between dual-vision methods and
 557 LVDT were 0.823 mm (Cam1-2), 0.507 mm (Cam1-3), and 0.949 mm (Cam2-3), respectively.
 558 The measurements from single vision and dual-vision matched well with the LVDT measurement.
 559 In addition, the measured displacements were found consistent for single-vision setting (Figure
 560 15) and dual-vision setting (Figure 16) on different floors with different extent of motion blur,
 561 validating the measurement robustness using multiple perspectives. Dual-vision improved the
 562 confidence in measurement compared with single-vision, although the accuracies of the two were
 563 similar in this study. When using single-vision, a strong assumption is required that allows only
 564 in-plane movement. The reason that both dual-vision and single-vision achieved similar accuracy
 565 in our study is that the excitation only caused in-plane vibration, meeting the required assumption
 566 for single-vision measurement. The implementation of augmentation strategies (GLDD and multi-
 567 vision) could address mild- and medium-level motion blur. However, when the motion blur is too

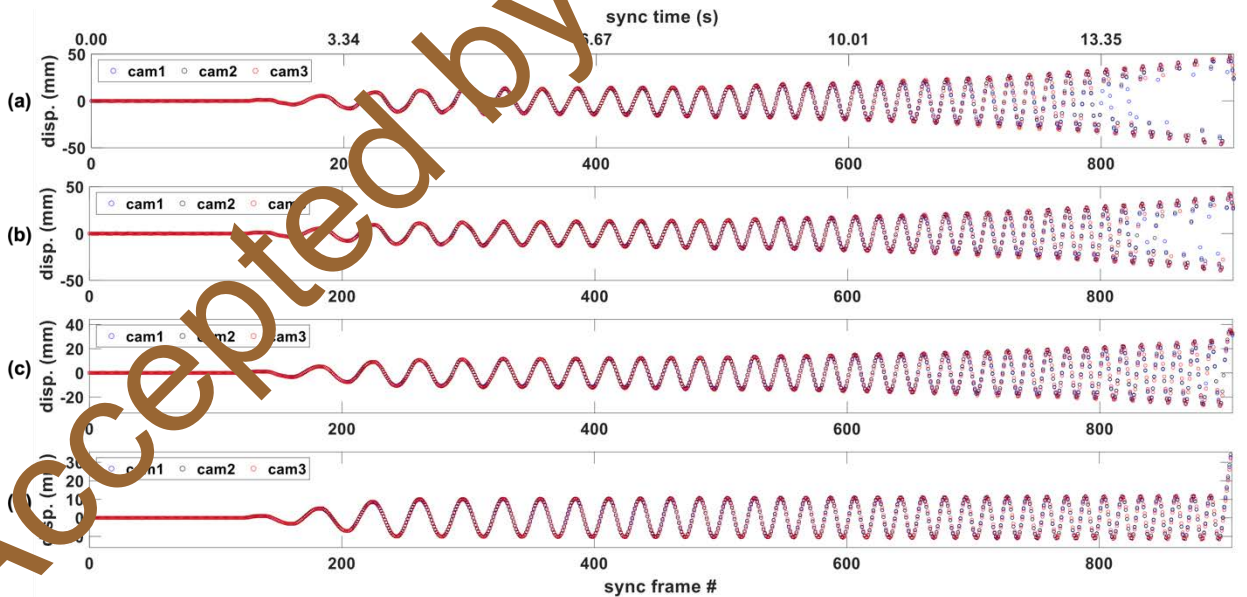
568 severe (e.g., vibration near natural frequencies), data fitting is needed to supplement and help
 569 interpolate/estimate the missing measurement.

570



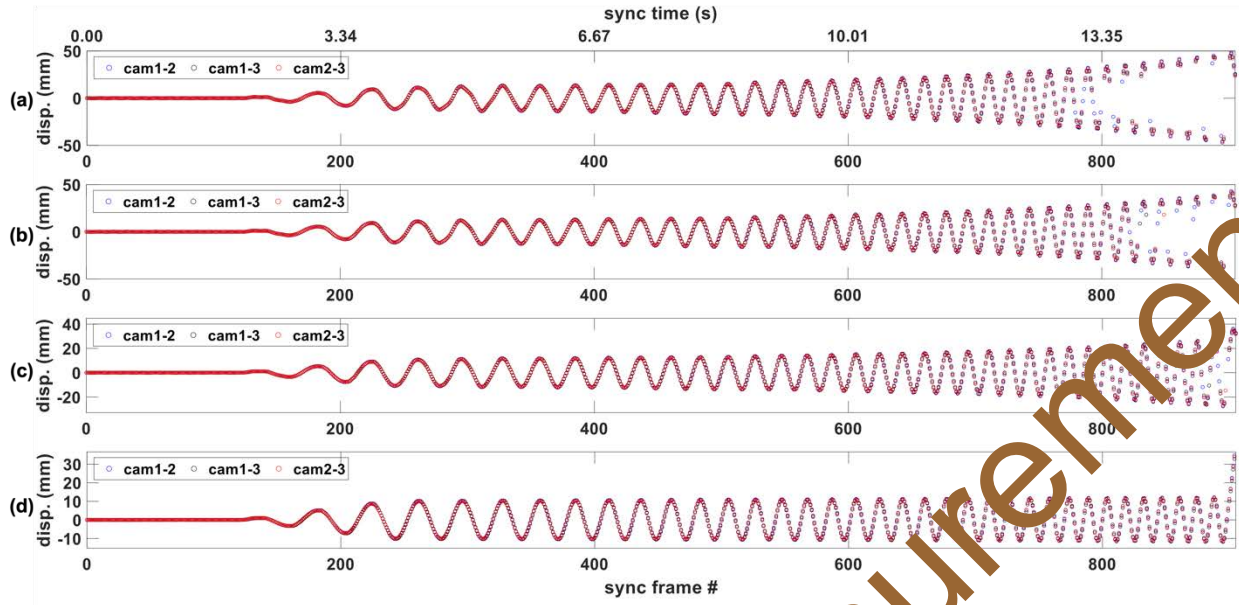
571

572 **Figure 14.** (a) Time histories of designed displacement input and LVDT measurement, (b) the differences
 573 between single-vision measurements and LVDT, and (c) the differences between dual-vision measurements
 574 and LVDT.



575

576 **Figure 15.** Displacement measured by single-vision method on (a) 3rd floor, (b) 2nd floor, (c) 1st floor, and
 577 (d) base.



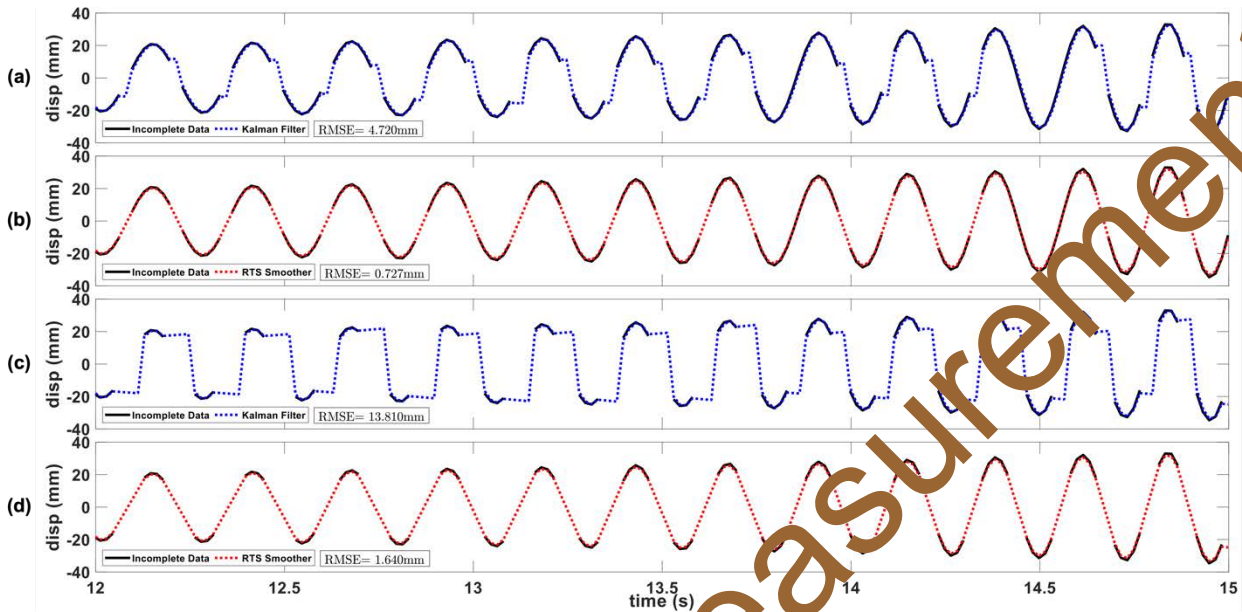
578

579 **Figure 16.** Displacement measured by dual-vision method on (a) 3rd floor, (b) 2nd floor, (c) 1st floor, and
 580 (d) base.

581 5.3. Data Fitting with Filter and Smoother

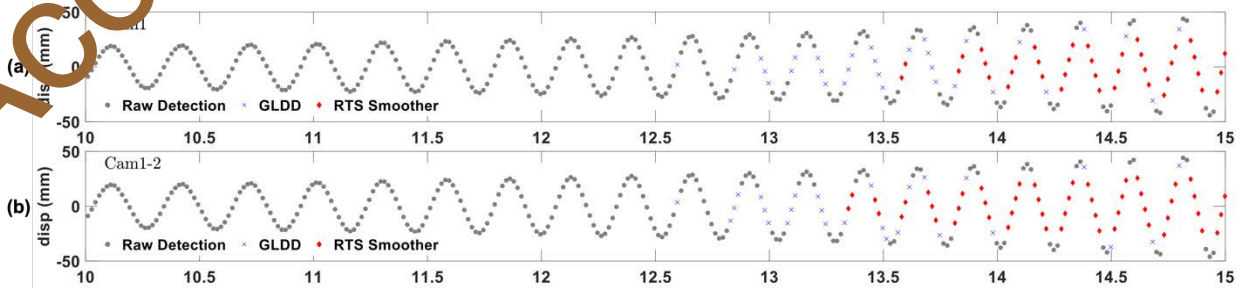
582 When motion blur was at mild or medium level in the experiment, the proposed GLDD
 583 method could restore some images for feature detections but could not resolve severe image blurs.
 584 Excessive motion blur was studied using filtering and smoothing methods to estimate the missed
 585 measurements. The virtual measurement on the 3rd floor of the FE model was used to evaluate the
 586 measurement estimation. Two virtual incomplete measurements between the two thresholds
 587 ($|d_1| \leq 10$ mm as shown in **Figure 17a-b**, and $|d_2| \leq 15$ mm as shown in **Figure 17c-b**) were
 588 masked with (k, i) within the time window of $12 s < t_k < 15 s$. The masked observations were treated
 589 as failed/missed observations ($(?_k)_i, (k, i) \in \mathcal{M}$). If the number of measurements for a single
 590 degree of freedom system is one, i can be dropped and the failed/missed observation can be
 591 presented as $(?_k), k \in \mathcal{M}$. Incomplete measurements were used to restore the unknown
 592 observation. The sampling time was $dt = 0.0167$ s to simulate vision-based measurement. The
 593 missed (virtual) measurements were within 12-15 s excluding the raw detections and the GLDD
 594 detection, to simulate the actual misdetection caused by the severe motion blur. In the KF and
 595 smoother setting, the initial system matrix was set as $\mathbf{A}_0 = [1, dt; 0, 1]$, the initial guess of the state
 596 was set as $\mathbf{x}_0 = [0, 0]^T$, and control was not considered in process equation. The covariance matrix

597 of the process \mathbf{Q}_0 was set as $[0.5, 0; 0, 0.5]$, and the dynamic model/transition matrix was set as
 598 $\mathbf{R}_k = 0.1 \text{ mm}^2$ based on the RMSE's in the evaluation of vision-based methods.



599
 600 **Figure 17.** Data fitting performance of Kalman filter and RTS smoother in two scenarios with (a-b)
 601 medium-level and (c-d) severe-level of (virtual) measurements from FE analysis.

602 **Figure 17** shows the performance of measurement fitting using KF and RTS smoother within
 603 the two incomplete time history data masked from 12-15 s. When there was a small number of
 604 missed measurements, e.g., 12.22% misses among the local time window of 12-15 s as 2.44%
 605 misses among the whole 0-15 s, KF still worked by neglecting the correction step (**Figure 17a**).
 606 However, when there were considerable number of missed measurements, e.g., 38.89% misses
 607 within the local time window of 12-15 s as 7.78% misses within the whole 0-15 s, the covariance
 608 matrix \mathbf{P}_k^- for state was enlarged without the necessary correction step. It was observed that the
 609 larger 13.810 mm RMSE occurred for KF estimation in severer blur case (**Figure 17c**) compared
 610 to 4.720 mm RMSE in mild blur case (**Figure 17a**).



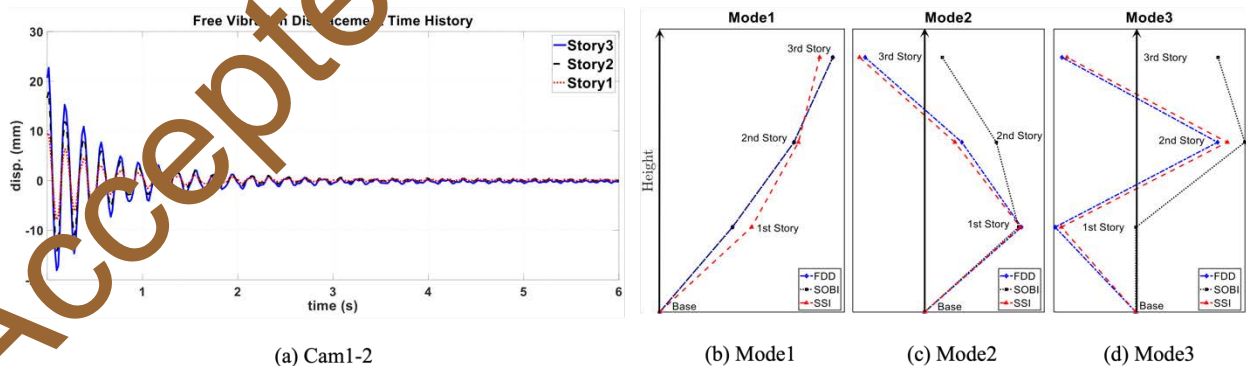
611

612 **Figure 18.** Data fitting performance of RTS smoother in experiment measurements from (a) single-vision
 613 (Cam1) and (b) dual-vision (Cam1-2).

614 The RTS data fitting took all available measurements into consideration, yielding better
 615 estimation performance (**Figure 17b-d**) with an improved RMSE of 0.727 mm and 1.640 mm for
 616 both mild and sever cases. The proposed RTS smoother-based fitting method was also
 617 implemented on the CV-based displacement measurements on the 3rd floor in the actual shake table
 618 test. The parameters of the RTS smoother for the experimental measurements were chosen as the
 619 same as the virtual one. As shown in **Figure 18**, the measurement from the raw detection and the
 620 measurement from the GLDD are denoted as gray dots and blue crosses, respectively. The
 621 estimation using the RTS smoother in single-vision case (**Figure 18a**) and dual vision case (**Figure**
 622 **18b**) are denoted as red diamond symbols. The data fitting results showed satisfactory estimation
 623 using RTS smoother.

624 5.4. Application of System Identification

625 After augmentation from the deblurring module and RTS smoother-based data fitting module,
 626 the measurement result can be used for system identification providing modal information for
 627 future applications (e.g., modal updating, structural damage identification). System identification
 628 can be based on structural displacements, such as free, forced, or ambient vibrations. To
 629 demonstrate the application of CV-based measurements in system identification, three output-only
 630 system identification methods of frequency domain decomposition (FDD) (Brincker et al., 2001),
 631 second-order blind identification (SOBI) (Belouchrani et al., 1997), and stochastic subspace
 632 identification (SSI) (Van Overchee & De Moor, 2012) were compared using both the virtual (from
 633 FE analysis) and experimental free vibration displacements.



634 (a) Cam1-2
 635 **Figure 19.** (a) Free-vibration displacements measured by dual-vision CV, and (b-d) the first three modal
 636 shapes from the system identification using FDD, SOBI, and SSI methods.

637 **Table 1** shows that the identified modal frequencies of the first three modes match well with
638 the FEM modal analysis from OpenSees. Modal assurance criterion (MAC) values for all three
639 modes are above 0.96 (except SOBI for mode-3 with 0.53) indicating a satisfactory performance
640 of FDD and SSI compared to SOBI. System identifications were also performed using the
641 experimental displacements (e.g., free vibrations in **Figure 19a**) measured by proposed method.
642 The identified modal frequencies (see **Table 1**) are very close with the 1st frequency identified as
643 5.15 Hz (FDD), 5.15 Hz (SOBI), and 5.14 Hz (SSI). The average differences between the identified
644 modal frequencies using FDD, SOBI, and SSI with the FE-based modal frequencies are 9.4%,
645 6.7%, and 4.3% for modes 1, 2, and 3, respectively. **Figure 19b-d** show the first three mode shapes
646 using the experimental measurements. It is found that the mode shapes between FDD and SSI are
647 close to each other, especially for mode-2 and mode-3. In general, the identification results were
648 consistent among the three methods using proposed multi-vision method.

649 6. Conclusion

650 The study proposed a multi-vision monitoring approach using low-cost cameras to measure
651 structural displacements in shake table tests with the augmentations from novel application of deep
652 learning-based image deblurring and Rauch-Tung-Striebel (RTS) Smoother. The proposed global-
653 local deblurring and detection (GLDD) module was able to restore clearer images for feature
654 detection, especially when dealing with mild level motion blur with average restoration rates of
655 92.1%. The restoration rates dropped to 50.6% for mild-level motion and further to 25.2% for
656 severe-level motion with the increasing severity of image blurs. Misdetections due to excessive
657 motion blur were estimated with filtering and smoothing-based methods using incomplete
658 measurements. RTS smoother is able to achieve a satisfactory data estimation (with a RMSE of
659 1.64 mm) outperforming Kalman filter (with a RMSE of 13.81 mm) in the scenario with severely
660 incomplete observations. RTS smoother helped accurately estimate missed measurement due to
661 severe blur, especially when the misdetections were consecutive as typical in shake table tests.
662 Implementation of GLDD module was tested in a shake table test of a three-story aluminum frame
663 and was validated with linear variable differential transformer measurement. Results show the
664 potential of the proposed approach in measuring dynamic displacement. The proposed multi-vision
665 and GLDD strategies can retrieve 75.0% measurements from previous misdetections (by just using
666 raw images from one single camera) and the data fitting module can complete the rest.

667 The main contribution of the study includes: (1) proposing a multi-vision displacement
668 measurement approach using low-cost cameras with novel deblurring module and RTS smoother-
669 based data fitting module to address the motion blur issue; (2) studying the effectiveness of the
670 modules in dealing with different levels of motion blur in shake table tests; (3) providing the
671 guidelines for using the proposed approach in shake table tests and the augmented displacement
672 that can be used in the further structural analysis. The proposed method can be employed in a range
673 of other applications (e.g., structural dynamics, finite element model updating) and be extended to
674 real-world applications, such as deflection measurement of bridge due to traffic loads, vibration
675 monitoring on high-rise buildings in earthquakes, and monitoring of relative displacement between
676 key structural members (e.g., inter-story, beam-column joints). One limitation of the study is that:
677 although the proposed multi-vision scheme and deblurring module is found to retrieve
678 misdetections due to mild and median motion blur, but it cannot restore images from excessive
679 motion blur, which is still a challenging issue for image processing. In addition, the effectiveness
680 of the proposed method using natural structural features under challenging environmental
681 conditions (e.g., poor illumination, occlusion of features), remains to be studied due to the limited
682 extent of this work. Future works will focus on studies of, such as effects of challenge conditions
683 (e.g., illumination, occlusion) in real application scenarios, faster algorithm on displacement
684 estimation, and error analysis of sensor deployment of the multi-vision system.

685

686 **Acknowledgment**

687 The authors thank Dr. Kevin Mackie from the University of Central Florida for providing the
688 shake table facilities. Help from Mr. Mostafa Iraniparast, Mr. Seyed Sina Shid-Moosavi, Ms. Lisa
689 Barra, and Ms. Aya Leon in the experiments is appreciated. This research is partially supported
690 by the Florida Department of Transportation (grant no.: BDV24 562-14 and BED26 562-4) and
691 the Florida Department of Environmental Protection William W. “Bill” Hinkley Center for Solid
692 and Hazardous Waste Management (award no.: AWD08952, project no.: P0184923).

693

694 **Author Contribution:**

695 The authors confirm contribution to the paper as follows: study conception and design: P. Sun;
696 data collection: M. Vasef; analysis and interpretation of results: P. Sun, M. Vasef; draft manuscript
697 preparation and revision: P. Sun, M. Vasef, L. Chen. All authors reviewed the results and approved
698 the final version of the manuscript.

699

700 **References**

701

702 Abdel-Aziz, Y. I., Karara, H. M., & Hauck, M. (2015). Direct linear transformation from comparator
703 coordinates into object space coordinates in close-range photogrammetry. *Photogrammetric
704 engineering & remote sensing*, *81*(2), 103-107.

705 Belouchrani, A., Abed-Meraim, K., Cardoso, J.-F., & Moulines, E. (1997). A blind source separation
706 technique using second-order statistics. *IEEE Transactions on signal processing*, *45*(2), 434-444.

707 Bezcioglu, M., Yigit, C. O., Karadeniz, B., Dindar, A. A., El-Mowafy, A., & Avci, Ö. (2023). Evaluation of
708 real-time variometric approach and real-time precise point positioning in monitoring dynamic
709 displacement based on high-rate (20 Hz) GPS Observations. *GPS Solutions*, *27*(1), 43.

710 Brincker, R., Zhang, L., & Andersen, P. (2001). Modal identification of output-only systems using
711 frequency domain decomposition. *Smart Materials and Structures*, *10*(3), 441.

712 Choi, H.-S., Cheung, J.-H., Kim, S.-H., & Ahn, J.-H. (2011). Structural dynamic displacement vision system
713 using digital image processing. *Ndt & E International*, *44*(7), 597-608.

714 Chu, A. (2005). CHAPTER 17 - Sensors for Mechanical Shock. In J. S. Wilson (Ed.), *Sensor Technology
715 Handbook* (pp. 457-480). Newnes. [https://doi.org/https://doi.org/10.1016/B978-075067729-
716 5/50057-4](https://doi.org/10.1016/B978-075067729-5/50057-4)

717 del Rey Castillo, E., Allen, T., Henry, R., Griffith, M., & Ingham, J. (2019). Digital image correlation (DIC)
718 for measurement of strains and displacements in coarse, low volume-fraction FRP composites
719 used in civil infrastructure. *Composite Structures*, *212*, 43-57.

720 Dong, C.-Z., Bas, S., & Catbas, F. N. (2020). A portable monitoring approach using cameras and computer
721 vision for bridge load rating in smart cities. *Journal of Civil Structural Health Monitoring*, *10*,
722 1001-1021.

723 Dong, C.-Z., & Catbas, F. N. (2021). A review of computer vision-based structural health monitoring at
724 local and global levels. *Structural Health Monitoring*, *20*(2), 692-743.

725 Gao, X., & Zhang, T. (2021). *Introduction to visual SLAM: from theory to practice*. Springer Nature.

726 Greenbaum, R. J., Smyth, A. W., & Chatzis, M. M. (2016). Monocular computer vision method for the
727 experimental study of three-dimensional rocking motion. *Journal of Engineering Mechanics*,
728 *142*(1), 04015062.

729 Guo, J., Jiao, J., Fujita, K., & Takewaki, I. (2019). Damage identification for frame structures using vision-
730 based measurement. *Engineering Structures*, *199*, 109634.

731 Han, K. a. W., Yunhe and Tian, Qi and Guo, Jianyuan and Xu, Chunjing and Xu, Chang. (2020). GhostNet:
732 More Features From Simpler Operations. *In Proceedings of the IEEE/CVF conference on computer
733 vision and pattern recognition 2020* (pp. 1580-1589).

734 Kogut, J. P., & Pilecká, T. (2020). Application of the terrestrial laser scanner in the monitoring of earth
735 structures. *Open Geosciences*, *12*(1), 503-517.

736 Krogius, M., Hagganmiller, A., & Olson, E. (2019). Flexible layouts for fiducial tags. 2019 IEEE/RSJ
737 International Conference on Intelligent Robots and Systems (IROS),

738 Kupyn, O., Budzai, V., Mykhailych, M., Mishkin, D., & Matas, J. (2018). DeblurGAN: Blind motion
739 deblurring using conditional adversarial networks. Proceedings of the IEEE conference on
740 computer vision and pattern recognition,

741 Kupyn, O. a. M., Tetiana and Wu, Junru and Wang, Zhangyang. (2019). DeblurGAN-v2: Deblurring
742 (Orders-of-Magnitude) Faster and Better. *Proceedings of the IEEE/CVF International Conference
743 on Computer Vision (ICCV), 2019*, pp. 8878-8887.

744 Liu, Y., Haridevan, A., Schofield, H., & Shan, J. (2022). Application of Ghost-DeblurGAN to Fiducial Marker
745 Detection. 2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS),

746 Ma, Y., Soatto, S., Košecká, J., & Sastry, S. (2004). *An invitation to 3-d vision: from images to geometric
747 models* (Vol. 26). Springer.

748 Ma, Z., Choi, J., & Sohn, H. (2022). Real - time structural displacement estimation by fusing
749 asynchronous acceleration and computer vision measurements. *Computer - Aided Civil and*
750 *Infrastructure Engineering*, 37(6), 688-703.

751 Mazzoni, S., McKenna, F., Scott, M. H., & Fenves, G. L. (2006). OpenSees command language manual.
752 *Pacific earthquake engineering research (PEER) center*, 264(1), 137-158.

753 Muralikrishnan, B. (2021). Performance evaluation of terrestrial laser scanners—A review. *Measurement*
754 *Science and Technology*, 32(7), 072001.

755 Olson, E. (2011). AprilTag: A robust and flexible visual fiducial system. 2011 IEEE international
756 conference on robotics and automation,

757 Ramakrishnan, S. a. P., Shubham and Gangopadhyay, Aalok and Raman, Shanmuganathan. (2017). Deep
758 Generative Filter for Motion Deblurring. *In Proceedings of the IEEE international conference on*
759 *computer vision workshops 2017 (pp. 2993-3000)*.

760 Rychlicki, M., Kasprzyk, Z., & Rosiński, A. (2020). Analysis of accuracy and reliability of different types of
761 GPS receivers. *Sensors*, 20(22), 6498.

762 Sandler, M. a. H., Andrew and Zhu, Menglong and Zhmoginov, Andrey and Chen, Liang-Chieh. (2018).
763 MobileNetV2: Inverted Residuals and Linear Bottlenecks. *In Proceedings of the IEEE conference*
764 *on computer vision and pattern recognition 2018 (pp. 4510-4520)*.

765 Särkkä, S. (2008). Unscented Rauch--Tung--Striebel smoother. *IEEE transactions on automatic control*,
766 53(3), 845-849.

767 Särkkä, S., & Svensson, L. (2023). *Bayesian filtering and smoothing* (Vol. 17). Cambridge university press.

768 Spencer Jr, B. F., Hoskere, V., & Narazaki, Y. (2019). Advances in computer vision-based civil
769 infrastructure inspection and monitoring. *Engineering*, 5(2), 199-222.

770 Sun, P., Bachilo, S. M., Lin, C. W., Nagarajaiah, S., & Wehman, R. B. (2019). Dual - layer nanotube -
771 based smart skin for enhanced noncontact strain sensing. *Structural Control and Health*
772 *Monitoring*, 26(1), e2279.

773 Sun, P., Draughon, G., Hou, R., & Lynch, J. P. (2022). Automated Human Use Mapping of Social
774 Infrastructure by Deep Learning Methods Applied to Smart City Camera Systems. *Journal of*
775 *Computing in Civil Engineering*, 36(1), 04022011.

776 Van Overschee, P., & De Moor, B. (2012). *Subspace identification for linear systems: Theory—*
777 *Implementation—Applications*. Springer Science & Business Media.

778 Wang, J., & Olson, E. (2016). AprilTag 2: Efficient and robust fiducial detection. 2016 IEEE/RSJ
779 International Conference on Intelligent Robots and Systems (IROS),

780 Welch, G., & Bishop, G. (1995). An introduction to the Kalman filter.

781 Xu, Y., Chen, D.-M., & Zhu, Y. (2019). Operational modal analysis using lifted continuously scanning
782 laser Doppler vibrometer measurements and its application to baseline-free structural damage
783 identification. *Journal of Vibration and Control*, 25(7), 1341-1364.

784 Yang, Y., Jung, Y. K., Horn, C., Park, G., Farrar, C., & Mascareñas, D. (2019). Estimation of full - field
785 dynamic strains from digital video measurements of output - only beam structures by video
786 motion processing and modal superposition. *Structural Control and Health Monitoring*, 26(10),
787 e2408.

788 Zeng, R., Wen, Y., Zhao, W., & Liu, Y.-J. (2020). View planning in robot active vision: A survey of systems,
789 algorithms, and applications. *Computational Visual Media*, 6, 225-245.

790 Zhang, D., Guo, J., Lei, X., & Zhu, C. (2016). A high-speed vision-based sensor for dynamic vibration
791 analysis using fast motion extraction algorithms. *Sensors*, 16(4), 572.

792 Zheng, W., Dan, D., Cheng, W., & Xia, Y. (2019). Real-time dynamic displacement monitoring with double
793 integration of acceleration based on recursive least squares method. *Measurement*, 141, 460-
794 471.

795 Zona, A. (2020). Vision-based vibration monitoring of structures and infrastructures: An overview of
796 recent applications. *Infrastructures*, 6(1), 4.

797

Accepted by Measurement