

Interactive Defect Quantification Through Extended Reality

Zaid Abbas Al-Sabbag^a, Chul Min Yeum^{a,*}, Sriram Narasimhan^b

^a*Department of Civil and Environmental Engineering, University of Waterloo, Canada*

^b*Department of Civil and Environmental Engineering, University of California, Los Angeles, USA*

Abstract

In this study, a new visual inspection method that can interactively detect and quantify structural defects using an *Extended Reality* (XR) device (headset) is proposed. The XR device, which is at the core of this method, supports an interactive environment using a holographic overlay of graphical information on the spatial environment and physical objects being inspected. By leveraging this capability, a novel XR-supported inspection pipeline, called *eXtended Reality-based Inspection and Visualization* (XRV), is developed. Key tasks supported by this method include detecting visual damage from sensory data acquired by the XR device, estimating its size, and visualizing (overlaid) information on the spatial environment. The crucial step of real-time interactive segmentation—detection and pixel-wise damage boundary refinement—is achieved using a feature Back-propagating Refinement Scheme (f-BRS) algorithm. Then, a ray-casting algorithm is applied to back-project the 2D image pixel coordinates of the damage region to their 3D world coordinates for damage area quantification in real-world (physical) units. Finally, the area information is overlaid and anchored to the scene containing damage for visualization and documentation. The performance of XRV is experimentally demonstrated by measuring surface structural damage of an in-service concrete bridge with less than 10% errors for two different test cases, and image processing latency of 2-3s (or 0.5s per seed point) from f-BRS. The proposed XRV pipeline underscores the advantages of real-time interaction between expert users and the XR device through immersive visualization so that a human-machine collaborative workflow can be established to obtain better inspection outcomes in terms of accuracy and robustness.

Keywords: visual inspection, extended reality, augmented reality, damage detection

1. Introduction

2. 1.1. Motivation

Large inventories of lifeline infrastructures supporting the water, energy, and transportation needs of our society are past or nearing the end of their intended service life and are deteriorating at an alarming rate. The potential failure of engineered structural components which comprise nearly all such infrastructures poses significantly increased risks to public safety [1], while simultaneously exacerbated

by the effects of climate change. New structural inspection guidelines in many jurisdictions have begun to adopt quantitative assessment criteria to evaluate their condition, replacing or augmenting qualitative indicators (e.g., estimating condition based on partial or faint visibility, or tapping areas of concrete with a hammer to measure delamination) [2]. For example, the *AASHTO Manual for Bridge Element Inspection* now requires inspectors to identify the types, locations, and sizes of damage on bridge elements to support more objective and quantitative estimates for deterioration [3]. Such a shift necessitates technology that can support an inspector in the field to achieve an accurate and objective assessment and satisfy new inspection requirements.

*Corresponding author at 200 University Ave. West, Waterloo, Ontario N2L 3G1, Canada

Email addresses: zaalsabb@uwaterloo.ca (Zaid Abbas Al-Sabbag), cmyeum@uwaterloo.ca (Chul Min Yeum), snarasim@ucla.edu (Sriram Narasimhan)

26 *Visual Inspection:* Fully manual visual assessment
27 remains the most widely used inspection method to-
28 day. This typically involves an inspector visually
29 observing the object and documenting various de-
30 fects observed (e.g., surface cracks, spalling, corro-
31 sion, etc.). Measurements of defects, in this case,
32 will typically involve the inspector physically ac-
33 cessing the area with a handheld measurement de-
34 vice. While this process can be simple and can pro-
35 duce the desired outcomes in many situations, es-
36 pecially when performed by an experienced inspec-
37 tor, often such methods are found to be inconsistent
38 (across multiple inspectors and inspections), cum-
39 bersome, costly, and unreliable (because of incorrect
40 measurements) [4]. New technology-driven stand-
41 alone inspection methods supported by advances in
42 computer vision, aerial and ground robots have tried
43 to address such challenges [5, 6, 7, 8]. While these
44 technologies provide the inspector with new means
45 to gather and process condition data from sensors
46 (cameras, lidars, etc.), they continue to rely on work-
47 flows that are mostly asynchronous and are often not
48 trustworthy. For example, under challenging situa-
49 tions (e.g., poor lighting, shadows, etc.), we cannot
50 expect to collect enough high-quality data to ensure
51 reliable outcomes, and inference from the automated
52 system alone may be unsatisfactory. Thus, the ex-
53 pertise and intelligence of the inspector should be in-
54 corporated into the process by adopting interactive
55 processes (e.g., checking data quality, refining pro-
56 cessed data). This paper aims to address this cru-
57 cial aspect by implementing superior data collection,
58 spatial mapping, and advanced computer vision al-
59 gorithms so that a human-machine collaborative in-
60 spection workflow can be established.

61 *XR technology:* XR broadly describes a range of
62 technologies that enable users to experience virtual
63 content in either a simulated or real environment. XR
64 encompasses *Augmented Reality* (AR), *Virtual Real-
65 ity* (VR), and *Mixed Reality* (MR) [9]. AR enables
66 visualization of virtual objects on video streams of
67 real scenes and is widely used in handheld digital de-
68 vices such as smartphones or tablets. However, inter-
69 activity between virtual objects, users, and the phys-
70 ical scene is limited because, in most cases, these
71 devices do not have a depth perspective, and thus
72 they cannot create a detailed spatial map of the real

73 scene. Users also have no way to directly interact
74 with 3D virtual objects, since most interactions in
75 AR are often mediated through 2D touch screens. On
76 the other hand, VR fully supports interactive envi-
77 ronments between virtual objects and users through
78 controllers or hand gestures, but this can only run
79 in the virtual world and there is no interaction with
80 the real-world scene. MR combines key elements of
81 AR and VR to let users see and interact with vir-
82 tual objects that are overlaid on the real environment
83 and can enable immersive visualization of informa-
84 tion as well as human-machine interaction [10, 11].
85 A key component of MR is the ability to digitize
86 the user’s environment by creating a spatial map of
87 the scene to facilitate interaction between real and
88 virtual objects and to replicate immersive visual ef-
89 fects such as occlusion or collision between virtual
90 objects and the real scene. Examples of MR dev-
91 ices that are currently available include Microsoft’s
92 HoloLens 1 (HL1), HoloLens 2 (HL2), and Magic
93 Leap One [12, 13, 14]. In this study, we used HL2
94 due to its intuitive user interface combined with ad-
95 vanced 3D spatial mapping capabilities of the phys-
96 ical environment which are crucial for the proposed
97 activities.

1.2. Literature Review

98 Significant literature exists on the topic of al-
99 leviating access through the use of aerial drones
100 and terrestrial mobile platforms equipped with vi-
101 sion cameras [15, 16, 17, 18, 19]. Similarly, auto-
102 mated detection of damaged areas from large vol-
103 umes of images using Deep Convolutional Neural
104 Networks (DCNNs) have also been studied exten-
105 sively in the context of computer vision applica-
106 tions to this problem [20, 21, 22, 23, 24, 25]. Si-
107 multaneous Localization and Mapping (SLAM) al-
108 gorithms [26, 27, 28, 29] enable inspectors to recon-
109 struct digital 3D point cloud models of the environ-
110 ment to localize and quantify damage in structural
111 elements [30, 31, 32]. However, a majority of these
112 works do not address the human-machine collabora-
113 tive aspect.

114 In this context, XR technologies have seen ap-
115 plications in medical, construction, manufac-
116 turing, education, entertainment, aerospace, and other
117 fields [33, 34, 35, 36]. However, their use in vision-

119 based infrastructure inspection and structural health
120 monitoring is relatively limited. In one of the
121 early works on the subject, AR was proposed for
122 the rapid assessment of earthquake-induced building
123 damage [37]. AR was also studied as a means to
124 annotate, organize and visualize images and meta-
125 data during inspections here [38]. Another research
126 work [39] demonstrated the application of HL1 for
127 bridge inspections by allowing users to measure
128 physical distances between any two points on the
129 structure. HL1 device has also been used to develop
130 a human-infrastructure interface for inspectors to
131 visualize and interact with real-time data received from
132 low-cost smart sensors mounted on structures [40].
133 In a related work, a framework for MR enabled
134 inspection system using AR/MR devices was recently
135 presented in [41]. The authors demonstrate that
136 their system can detect and segment damage using
137 a semi-supervised SegNet-based DCNN, where
138 users can modify the prediction threshold in real-
139 time to improve segmentation results. However, the
140 system was implemented on an AR device without
141 spatial mapping or 3D sensing capabilities, limiting
142 quantitative evaluation. All of these studies have
143 addressed the collaborative aspects of XR technol-
144 ogy. However, they did not harness the spatial map-
145 ping capabilities provided by XR devices—through
146 SLAM—to localize damage in the real scene and
147 quantitatively evaluate inspection regions with high
148 accuracy, which can limit the potential usage of XR
149 devices in the actual inspection process.

150 1.3. Objectives

151 In this study, a novel real-time visual interac-
152 tive inspection system is proposed, which can de-
153 tect and measure the area of surface damage enabled
154 through XR (HL2), called *eXtended Reality-based*
155 *Inspection and Visualization* (XRIV). XRIV allows
156 inspectors to capture images of a structure, detect vi-
157 sual damage on a surface, estimate its area, and then
158 visualize them as holographic objects overlaid on the
159 scene. Data collection and analysis are integrated
160 into a single synchronous on-site process. A DCNN-
161 based interactive segmentation algorithm called fea-
162 ture Back-propagating Refinement Scheme (f-BRS)
163 is utilized to segment the damage region on an im-
164 age and refine the results interactively through user

165 feedback. Once the damage region is segmented in-
166 teractively, the 3D world coordinates of this region’s
167 boundary are estimated through back-projection us-
168 ing the pose of the camera and the spatial map cre-
169 ated by the device, followed by area estimation.

170 The main novelty of this work is to show how
171 discrete hardware and software components within
172 XR, and advanced computer vision algorithms such
173 as interactive segmentation, are customized and in-
174 tegrated into a single, end-to-end workflow to create
175 an interactive visual inspection system by involving
176 the expert user to obtain and refine real-time quan-
177 titative damage measurements. Integration toward
178 enhancing interactivity between users and the sys-
179 tem results in more accurate and efficient visual in-
180 spection in the field. Also, a novel procedure was
181 devised and deployed to spatially match 2D image
182 pixels to the 3D scene using a ray-casting algorithm,
183 followed by estimating the physical size of the dam-
184 age. We experimentally demonstrated an end-to-end
185 visual inspection system by leveraging the immer-
186 sive and interactive spatial mapping capabilities pro-
187 vided by XR and data-driven reasoning achieved by
188 computer vision. We successfully established a fully
189 collaborative environment for structural inspections
190 and measurements of visual defects. Fig. 1 captures
191 XRIV in action, where it creates an interactive envi-
192 ronment for the user (through hand gestures) to mea-
193 sure the area of surface concrete spalling damage.
194 Fig. 1a shows the user’s point-of-view when using
195 XRIV implemented on HL2, while Fig. 1b shows the
196 same scene from a third-person view through a con-
197 nected Android smartphone.

198 This paper is organized as follows. First, Sec-
199 tion 2 covers the technical background of the key
200 XRIV components: interactive segmentation, spa-
201 tial computation, and area computation. Section 3
202 then outlines how these components are integrated
203 as part of the full system implementation and how
204 computational tasks are distributed between the XR
205 device and the server to achieve near real-time anal-
206 ysis. Section 4 shows the results of the several ex-
207 perimental tests that were conducted to validate the ac-
208 curacy of the f-BRS algorithm for spalling segmen-
209 tation, the spatial measurement accuracy of HL2,
210 and to demonstrate the performance of XRIV in the
211 field. Section 5 summarizes the key conclusions of



Figure 1: The proposed XRIV system for damage area measurement enabled through the HoloLens 2. Image (a) was captured using the HoloLens 2, while image (b) was captured using Spectator View [42] from a connected Android smartphone that views the same scene as HoloLens 2 through a two-way interface. Note that the smartphone is only used to capture image (b) and is not required for XRIV.

212 the study.

213 2. Technical Background

214 2.1. Interactive Segmentation

215 Most vision-based methods for detecting and
216 quantifying damage have utilized automated DCNN-
217 based methods to classify images when damage may
218 be present [20, 21, 25], or to perform automatic
219 pixel-wise segmentation of the damage region [23,
220 24, 43]. Although the accuracy of these methods
221 has dramatically improved over the years, there are
222 still some challenges. For example, automated algo-
223 rithms do not provide the user with a direct way to
224 refine incorrect segmentation results. Any segmen-
225 tation algorithm, no matter how accurate it is, will
226 inevitably encounter ambiguous situations where its
227 performance degrades. Even in routine settings, the
228 boundary of damage is not always clear and varies
229 depending on how damage is defined (e.g., concrete
230 flaking vs. spalling). In such cases, human judgment
231 and expertise can significantly improve the outcome
232 and leverage input from users to remove false detec-
233 tion and to enhance the performance of segmenta-
234 tion algorithms. Therefore, XRIV leverages a state-
235 of-the-art interactive segmentation algorithm, called
236 f-BRS [44], to detect and segment damage in im-
237 ages captured using an XR headset. The proposed
238 XRIV framework utilizes the holographic user inter-
239 face and f-BRS to enable inspectors to be involved

240 in the detection and segmentation process by being
241 able to correct inaccurate results with minimal user
242 interaction.

243 DCNN-based interactive segmentation algorithms
244 fall under the umbrella of semi-supervised learning,
245 which are machine learning algorithms that use a
246 combination of labelled and unlabelled data for train-
247 ing or to perform predictions [45]. With prior train-
248 ing on unlabelled data, and with some limited la-
249 belling provided by the user, interactive segmenta-
250 tion can segment object classes without being trained
251 on those classes. Users first select multiple points
252 on the image, called seed points, inside (positive)
253 and outside (negative) a target region of interest [46].
254 Seed points provide the DCNN with initial con-
255 ditions to separate the foreground (damage region in
256 this case) from the background. The seed points are
257 used to produce positive and negative distance maps
258 that encode information about where the points were
259 selected and placed on the image. These maps are
260 then concatenated to the image array and fed to the
261 DCNN as multi-layered inputs. The output of the
262 DCNN is a prediction map with confidence values
263 ranging from 0 to 1 for every pixel in the input im-
264 age. The confidence values represent the probability
265 of how likely each of the pixels is included in the
266 foreground, such as structural damage. The predic-
267 tion map is then converted to a binary segmenta-
268 tion mask by thresholding the confidence values to be-
269 come either 0 or 1.

270 An issue with DCNN-based interactive segmenta-
271 tion is that user-provided seed points are not always
272 entirely consistent with output prediction maps. This
273 is because a simple feed-forward DCNN model uses
274 back-propagation to train the network to minimize
275 overall segmentation error but does not constrain the
276 prediction map to be consistent with seed points dur-
277 ing the training phase [46]. For example, when the
278 user selects a positive seed point at a certain location
279 in the foreground of the image, the value of the pre-
280 diction map at that location may be slightly less than
281 1, even though the user is confident that this location
282 is part of the foreground. These small inconsistencies
283 accumulate and increase the overall error of the seg-
284 mentation, resulting in inaccurate boundary detec-
285 tion regardless of the number of seed points. A Back-
286 propagating Refinement Scheme (BRS) solves this
287 issue by running a backward pass through the net-
288 work for each selected seed point to update the input
289 distance maps and enforce an additional constraint
290 to make the output prediction map consistent with
291 the selected seed points' locations [47]. BRS pro-
292 duces prediction maps that are consistent with user-
293 selected seed points, producing higher overall accu-
294 racy for segmentation masks with fewer seed points.
295 However, BRS is a computationally expensive algo-
296 rithm and may produce significant delay because it
297 requires running a backward pass through the DCNN
298 for every selected seed point.

299 To address this issue, f-BRS builds upon the
300 concept of back-propagation refinement but utilizes
301 backward passes to modify the values of feature
302 maps of selected intermediate layers instead of mod-
303 ifying the inputs of the network. This means that
304 f-BRS requires running backward passes from the
305 output of the DCNN to the intermediate layers while
306 skipping the initial layers for each pass. Thus, f-BRS
307 can achieve the same level of accuracy as BRS but re-
308 duce computational cost per seed point because each
309 backward pass is shorter. This makes f-BRS ideal for
310 real-time interactive segmentation applications.

311 The segmentation model in Xriv builds on a pre-
312 trained f-BRS model that utilizes a ResNet-34 back-
313 bone pretrained on ImageNet, and a DeepLabV3+
314 decoder [48] containing the two intermediate convo-
315 lutions layers that were trained on the benchmark Se-
316 mantic Boundaries Dataset (SBD) [49]. Three f-BRS

317 configurations are considered in this study (labeled
318 A, B, and C), which differ based on the intermediate
319 layers to which the back-propagation refinement is
320 applied in the network. f-BRS-A applies the refine-
321 ment scheme to the output from the ResNet back-
322 bone, f-BRS-B applies it inside the DeepLabV3+
323 decoder but before the first convolution layer, and
324 f-BRS-C applies it before the second convolution
325 layer. The performance of each configuration was
326 experimentally compared in Section 4.1.

327 In addition to the variation in the network configura-
328 tion, a key difference between the proposed and the
329 original implementation [44] is in how seed points
330 are used for segmentation. The original implemen-
331 tation was designed to be used on a static image to al-
332 low users to select one seed at a time by successively
333 placing points on incorrectly segmented regions to
334 improve the result. However, direct implementa-
335 tion of this process is not suitable for XR applica-
336 tions because processing each seed point is computa-
337 tionally inefficient, and the user would need to wait for
338 the previous seed point to be processed before se-
339 lecting the next seed point, which can become time-
340 consuming if many seed points need to be selected
341 to obtain a reasonably accurate result. Instead, the
342 f-BRS algorithm was customized for real-time im-
343 plementation in Xriv to allow users to select mul-
344 tiple seed points and then capture an image for con-
345 ducting segmentation at once, instead of processing
346 each seed point individually. This can speed up the
347 process considerably without an associated penalty
348 in the accuracy because image segmentation would
349 only need to be run once unless the user decides to
350 modify the segmentation result.

351 2.2. Spatial Computation

352 A key component of Xriv is spatial computation,
353 which is the process of computing how points in the
354 3D world coordinate system are projected to the 2D
355 pixel coordinate system of the built-in camera and
356 vice versa, as shown in Fig. 2. The XR headset's
357 built-in camera follows a pinhole camera geometry
358 model and stores a 3×4 camera projection matrix (P)
359 for each image or video. P contains camera intrinsics
360 and the pose of the camera (obtained from SLAM in
361 the device) synchronized to the time of image cap-
362 ture. Since the camera parameters are known (pro-

Algorithm 1 Ray-casting

Input

C : Ray origin point (Camera Optical Center)
 V : Ray direction unit vector (P^+x)
 M : Set of m triangles in spatial mesh

Output

X : Intersection point between ray and spatial mesh

```

1:  $\lambda_{min} \leftarrow \text{empty}$                                 ▷ Initialize  $\lambda_{min}$  as a placeholder value
2: for triangle  $i = 1, 2, \dots, m$  in  $M$  do      ▷ Iterate over triangle  $i$  joined by points  $A_i, B_i, C_i$  with normal  $N_i$ 
3:    $\lambda_i \leftarrow \frac{(C - A_i) \cdot N_i}{V \cdot N_i}$           ▷ Find  $\lambda_i$  using ray-plane intersection equation
4:   if  $\lambda_i > 0$  then                                ▷ Check if the intersection is in front of the camera
5:      $X_i \leftarrow C + \lambda_i V$                       ▷ Assign candidate intersection point  $X_i$ 
6:      $v_0 \leftarrow C_i - A_i$ 
7:      $v_1 \leftarrow B_i - A_i$ 
8:      $v_2 \leftarrow X_i - A_i$ 
9:      $\begin{bmatrix} \alpha \\ \beta \end{bmatrix} \leftarrow \begin{bmatrix} v_0 \cdot v_0 & v_1 \cdot v_0 \\ v_0 \cdot v_1 & v_1 \cdot v_1 \end{bmatrix}^{-1} \begin{bmatrix} v_2 \cdot v_0 \\ v_2 \cdot v_1 \end{bmatrix}$     ▷ Find barycentric coordinates  $\alpha, \beta$  of  $X_i$  on triangle surface
10:
11:   if  $0 \leq \alpha \leq 1$  and  $0 \leq \beta \leq 1$  and  $0 \leq \alpha + \beta \leq 1$  then    ▷ Check if  $X_i$  is inside triangle
12:     if  $\lambda_i < \lambda_{min}$  then                                ▷ Assign  $\lambda_i$  to  $\lambda_{min}$  if it is the lowest calculated value yet
13:        $\lambda_{min} \leftarrow \lambda_i$ 
14:     end if
15:   end if
16: end if
17: end for
18:  $X \leftarrow C + \lambda_{min} V$                                 ▷ Calculate intersection point closest to origin point using  $\lambda_{min}$ 

```

363 vided by the manufacturer), the estimation is limited
364 to the 6 DOF pose of the camera associated with an
365 image, which includes 3D translation and rotation
366 parameters. P can then be constructed using a known

367 3×3 camera intrinsic matrix (K), a 3×3 camera rotation
368 matrix (R), and a 3×1 camera translation vector
369 (t), as presented in Eq. (1).

$$P = K \begin{bmatrix} R & | & t \end{bmatrix} \quad (1)$$

The 2D pixel coordinates (x_s) of points along the boundary of the segmented damage region on the image plane and the corresponding 3D points (X_s) in the world coordinates are related through the projective transformation in Eq. (2):

$$x_s = PX_s \quad (2)$$

Here, the 2D pixel coordinates and 3D world coordinates are represented in the homogeneous coordinate system as 3×1 vectors $(u, v, 1)^T$ and 4×1 vectors $(x, y, z, 1)^T$, respectively. This transformation is also used to project the world coordinates of user-selected

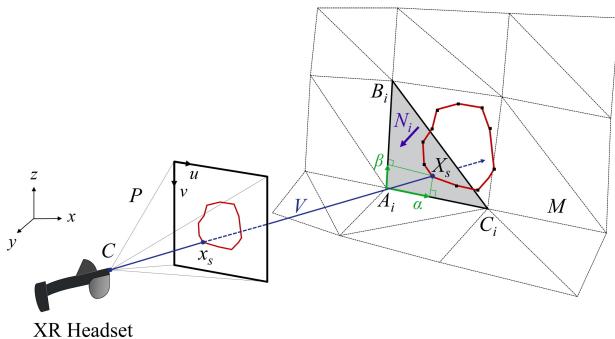


Figure 2: Pinhole camera geometry: Eq. (2) describes the projection of 3D world points to image points. A ray-casting algorithm in Algorithm 1 is used to back-project 2D image pixel coordinates of damage boundary to spatial mesh in 3D.

380 seed points (X_p) onto the image plane to obtain their
 381 image pixel coordinates (x_p).
 382

The projective transformation in Eq. (2) allows X_s
 383 to be directly mapped to pixel coordinates on the im-
 384 age plane. However, since P is a non-square non-
 385 singular matrix, inverting this transformation yields
 386 non-unique solutions for X_s that correspond to x_s .
 387 A unique solution for X_s is obtained through back-
 388 projection, where image pixels are extended into the
 389 3D world coordinates as vectors to determine associa-
 390 tions, a process known as ray-casting. This process
 391 is used to anchor points on the spatial mesh of the
 392 3D environment and to obtain the physical scale of
 393 the damage region for area computations.

The details of the ray-casting algorithm imple-
 394 mented in this paper are described in Algorithm 1.
 395 With reference to Fig. 2, the process of ray-casting
 396 locates the intersection point of the ray extended
 397 from the camera optical center (C) from a point
 398 on the image plane with the spatial mesh (M) con-
 399 structed by the XR device. Suppose that the ray has
 400 an origin point at C along a certain ray direction vec-
 401 tor (V) specified by an image point (x) and with an
 402 unknown distance parameter (λ), then the points on
 403 the ray (X) can be represented as a parametric line
 404 equation in Eq. (3) [50]:
 405

$$X(\lambda) = C + \lambda V \text{ where } V = P^+ x \quad (3)$$

406 where P^+ is the pseudo-inverse of P . M consists of
 407 a set of triangles in the world coordinates that are
 408 joined along their edges, as seen in Fig. 2.
 409

The ray-casting algorithm takes C , V , and M as
 410 inputs and iterates over every triangle in M to deter-
 411 mine an intersection point (X_i) between the ray and
 412 each triangle using a ray-plane intersection equation.
 413 Then, each X_i is checked to determine if it lies inside
 414 the triangle using a point-in-triangle test [51]. The
 415 point X_i is first converted to 2D barycentric coordi-
 416 nates (α, β) that indicate where the point is relative
 417 to the triangle's vertices, denoted as A_i, B_i and C_i . If
 418 α, β , and their sum value are between 0 and 1, then
 419 X_i is determined to be inside the triangle, which is
 420 the intersection point. However, a ray can intersect
 421 with M at multiple points since the real scene of M
 422 could be complex. In order to select a single out-
 423 put intersection point (X) while removing all other

points that are either occluded or behind the camera,
 424 the X_i associated with the minimum λ (λ_{min}) greater
 425 than 0 is retained as the output. Therefore, every 2D
 426 image point (x) can be mapped to a unique 3D world
 427 coordinates point (X) on M using ray-casting.
 428

Algorithm 1 is used to determine the world coor-
 429 dinates of each point along the segmented damage
 430 boundary by setting $x = x_s$ and $X = X_s$ for every im-
 431 age point of the boundary. Ray-casting is also used to
 432 determine the world coordinates of seed points (X_p)
 433 selected by the user through hand gestures using the
 434 XR interface. In this case, V in Algorithm 1 is de-
 435 termined through hand gesture, while C denotes the
 436 world coordinates of the wrist where the ray is as-
 437 sumed to originate.
 438

2.3. Area Computation

Once the 3D world coordinates of the damage re-
 440 gion boundary have been determined, the physical
 441 area of the bounded region is computed. To compute
 442 the areas in world coordinates, the general assump-
 443 tion is that the boundary lies on a planar surface with
 444 a single normal vector (N_s). However, due to the
 445 non-flatness of the mesh surface being analyzed as
 446 well as errors that are present in spatial mesh, the
 447 boundary points are not always placed on the same
 448 plane.
 449

To overcome this issue, a best-fit reference plane
 450 is first determined by minimizing the least-squares
 451 error of the boundary points. The world coordinates
 452 of the boundary (X_s) are projected onto the reference
 453 plane with normal vector N_s to obtain a set of 2D
 454 planar coordinates (y_s) using a projective geometric
 455 transformation (T) by making one of three dimen-
 456 sions as zero. The points are first shifted to the cen-
 457 troid of the region by subtracting the average of the
 458 boundary points (\bar{X}_s) before pre-multiplying them by
 459 T in Eq. (4). The transformation is constructed by
 460 concatenating any two ortho-normal column vectors,
 461 denoted as n_u and n_v , that lie on the reference plane
 462 and are perpendicular to N_s .
 463

$$y_s = T(X_s - \bar{X}_s) \text{ where } T = \begin{bmatrix} n_u & n_v \end{bmatrix}^T \quad (4)$$

The set of 2D planar boundary coordinates, de-
 464 noted as y_s , describe a closed polygon joined by k
 465

number of points. The area of this polygon (A_s) is then calculated using the Shoelace formula [52] in Eq. 5, which is used to subdivide the polygon into k triangles and add up the areas of all of them. For every point (u_i, v_i) connected to the next point (u_{i+1}, v_{i+1}) in y_s , the area of a triangle connecting both points to the origin is calculated using a matrix determinant. A_s is then calculated by summing up the areas of all the triangles,

$$A_s = \frac{1}{2} \sum_{i=1}^k \begin{vmatrix} u_i & u_{i+1} \\ v_i & v_{i+1} \end{vmatrix} \quad (5)$$

3. System Implementation

The implementation schema of XRIV includes component subsystems where tasks are apportioned between the XR device and a computation server to allow for efficient use of available computing resources to enable real-time analysis. The user input/output interaction systems enabled through the XR interface allow users to interact with holographic objects such as the UI menu (shown in Fig. 1) using hand gestures and to visualize the segmented damage region and its area information through holographic overlays. Spatial computation using SLAM is executed on the XR device and the remaining heavy computational tasks including DCNN-based interactive segmentation and area computation are offloaded to a remote computational server. Here, the computational server can include, for example, cloud computing or a physical device such as a smartphone or a mobile computer. The detailed implementation pipeline of XRIV is described in Fig. 3. The interactive components of XRIV, namely locating damaged regions, the selection of the seed points, and interactively refining segmentation results to quantify damaged regions are illustrated in Fig. 4.

One of the key inputs for interaction is actualized in the component *Input User Interaction* in Fig. 3, namely selecting input seed points to detect damage on structures using interactive segmentation and to refine their results. First, the user selects positive and negative seed points from inside and outside the target region, respectively (Step 1 in Fig. 4a), to assist the DCNN with separating the foreground (damage

region) from the background. Hand gesture recognition in the XR device allows the user to place the virtual seed points by interacting with the XR device's UI features. Upon receiving this input, the XR device then calculates the 3D world coordinates (X_p) of the seed points using a ray-casting algorithm, denoted *rayCast*, to anchor the points to the spatial mesh (M) (Algorithm 1). Then, the user captures an image (I) using the XR device's built-in camera that must show the full target damage region as well as all the seed points selected, initiating Step 2 in Fig. 4b. The 2D pixel coordinates (x_p) in I corresponding to X_p are then obtained by multiplying X_p by the corresponding projection matrix (P) (see Eq. 2 in Section 2.2).

Next, the XR device sends I and x_p to the computational server to segment the damage region using the f-BRS segmentation algorithm, denoted as *Segment* in Fig. 3 (Section 2.1). *Segment* takes x_p as input seed points to segment I and produces a binary mask (B) that indicates where the damage is present in I . A topological structural analysis algorithm, denoted as *findContours*, is used to extract contours from B to obtain the 2D pixel coordinates of the boundary (x_s) of the segmented region [53]. The 2D pixel coordinates, x_s , of the damage boundary extracted from I are then sent back to the XR device. Note that the computational offloading of *Segment* and *findContours* through the remote computation server in XRIV enables computationally expensive tasks to be undertaken in a remote machine rather than on a wearable device with limited computing and power resources.

The pseudo-inverse P^+ calculation and the back-projection of x_s as a set of continuous rays originating from the camera optical center (C) of I are executed on the XR device (denoted by *rayCast* in Fig. 3) to obtain a set of 3D world coordinates of the damage region boundary (X_s) corresponding to x_s . The XR interface then displays a holographic overlay of the target damage region and anchors it to M (red colored region in Fig. 4b).

In *User Output Interaction* in Fig. 3, the user evaluates whether the target damage region is properly segmented based on the holographic overlay. In some cases, damage regions might be over-or under-estimated such as the one shown in Fig. 4b. This component allows the user to interactively improve

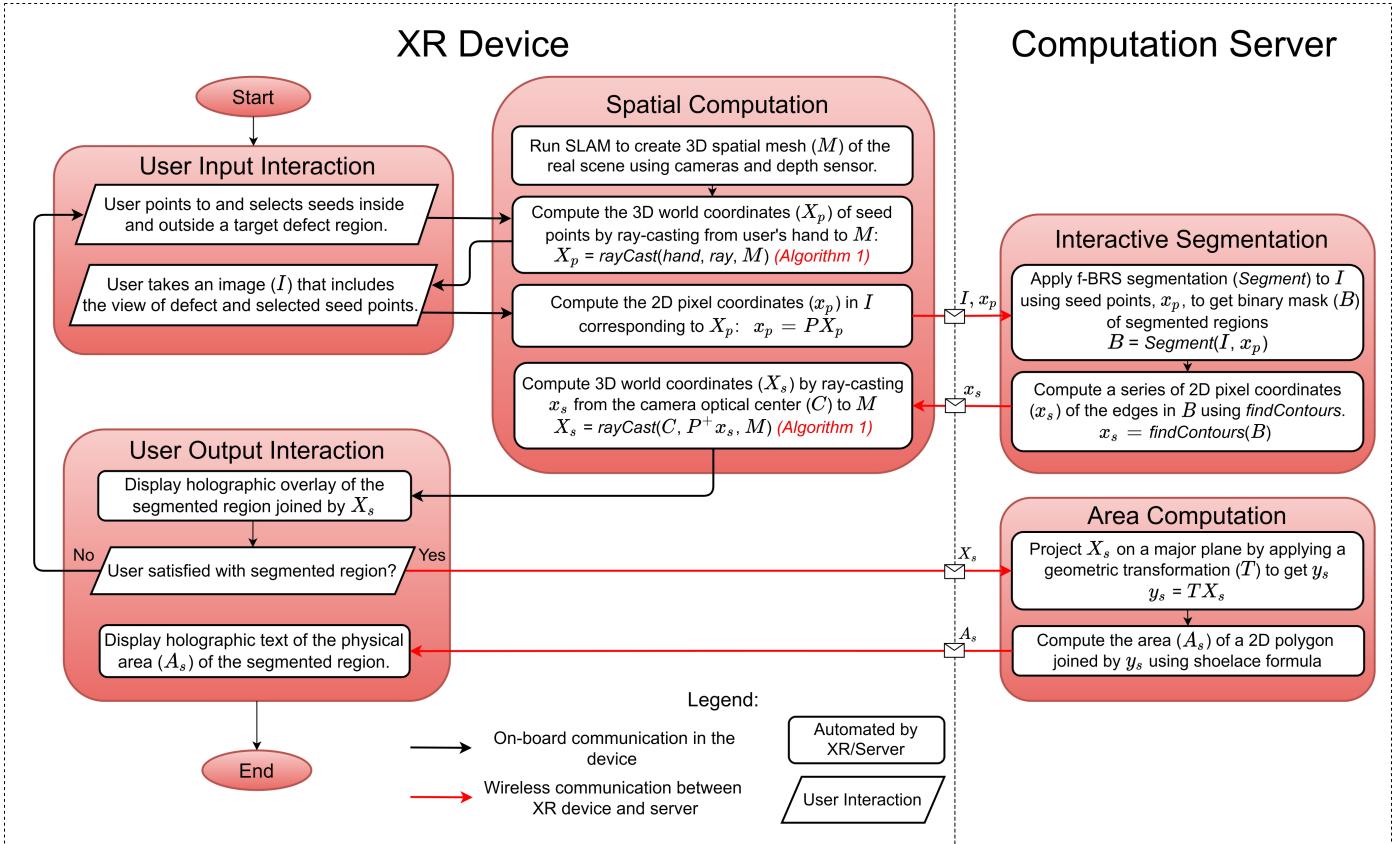


Figure 3: Real-time Xriv pipeline using XR device and computational server. Wireless communication is through HTTP requests.

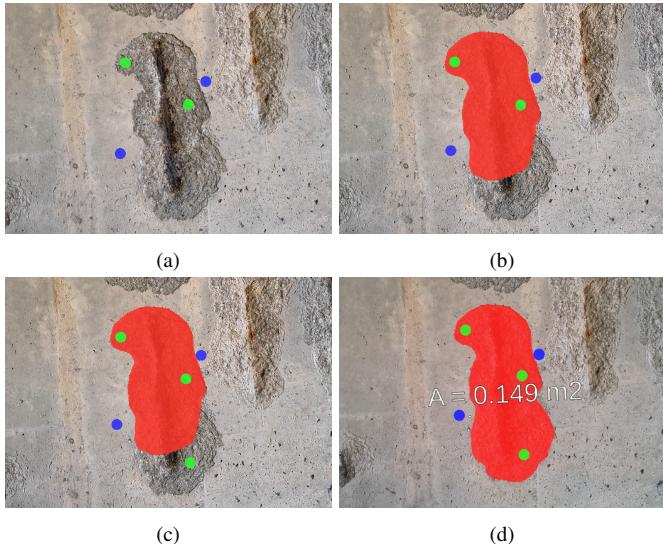


Figure 4: XR interface for interactive damage segmentation, demonstrated on the HoloLens 2: (a) Step 1: select positive and negative seed points inside and outside damage region, (b) Step 2: apply interactive segmentation algorithm to the captured image, (c) Step 3: add seed points to refine segmentation results, and (d) Step 4: calculate the area of the refined segmented region.

the quality of segmentation, by adding more positive or negative seed points and then repeating the segmentation process (Step 3 in Fig. 4c). After obtaining a satisfactory segmented damage region, the user presses a button on the holographic menu (shown in Fig. 1) to obtain the physical area (A_s), which is mathematically represented as a closed polygon connected by X_s (see Section 2.3). Finally, a holographic text of the physical area of the target damage region is displayed in the XR interface and anchored to the target damage region (see Step 4 in Fig. 4d).

In addition to real-time analysis capability, Xriv also supports offline processing of these tasks. Such offline processing may be desirable to online processing, e.g., when several damage regions need to be inspected, or when a field inspection should be completed within a short period and the inspector prefers to analyze the data and monitor its quality at a later date. Although such an asynchronous offline mode does not facilitate in-situ real-time analysis, it allows for more flexibility in the selection of segmentation algorithms or the method of segmen-

576 tation. For example, the user can provide one seed
 577 point at a time while checking the quality of seg-
 578 mentation, rather than providing a set of seed points
 579 at once (explained in Section 2.1), or can implement
 580 other detection or segmentation algorithms.

581 4. Experimental Results

582 4.1. f-BRS Segmentation Evaluation

583 To evaluate the performance of f-BRS for dam-
 584 age (spalling) segmentation, different implemen-
 585 tations were tested by varying the seed points provided
 586 by the user. A total of 91 images of spalling damage
 587 were collected from in-service bridges in Southern
 588 Ontario and the boundary of the spalling damage was
 589 manually annotated. Three implementations of f-
 590 BRS and the original BRS were first compared. The
 591 original f-BRS study uses the evaluation metric of
 592 Intersection-over-Union (IoU) to measure segmenta-
 593 tion error between the predicted and the ground-truth
 594 regions [44]. However, since the area measurement
 595 accuracy is a major concern rather than IoU in this
 596 paper, the evaluation metric chosen here is the area
 597 measurement accuracy, which is defined as the pre-
 598 dicted area (A_P) as a percentage of the ground-truth
 599 area (A_G).

600 The first step was to randomly select a pair of
 601 positive and negative seed points which were placed
 602 inside and outside the ground-truth region, respec-
 603 tively. With these seed points, a predicted mask was
 604 generated and compared with the ground-truth re-
 605 gion to determine the area measurement accuracy.

In the next iteration, the set of positive and negative
 606 seed points were randomly selected in locations that
 607 were incorrectly segmented to generate a more accu-
 608 rate segmentation mask. This process was repeated
 609 for up to 12 pairs of seed points to incrementally re-
 610 fine the segmentation results. Fig. 5 presents sample
 611 segmentation results when the area measurement ac-
 612 curacy using f-BRS-B is over 80%. The positive (in-
 613 side spalling region) and negative (outside spalling
 614 region) seed points used for the segmentation of each
 615 image are displayed over the segmentation masks.

The results of the experiment, presented in Fig. 6
 617 show the average area measurement accuracy of all
 618 damage in images, and are computed depending on
 619 the number of seed points pairs. Overall, all im-
 620 plementations of f-BRS, as well as BRS, can seg-
 621 ment spalling-type damage with 80% accuracy on
 622 average if five or more seed point pairs (10 seed
 623 points in total) are selected, which is satisfactory for
 624 vision-based inspection applications [2, 3]. Over-
 625 all, f-BRS-B achieves slightly better accuracy among
 626 the f-BRS implementations and approaches a simi-
 627 lar level of accuracy as BRS. Note that f-BRS was
 628 designed to improve the speed of the segmentation
 629 in BRS without significantly sacrificing accuracy.
 630 Therefore, the accuracy of BRS becomes an upper
 631 limit for f-BRS. For this experiment—performed on
 632 a computer equipped with a GPU: NVIDIA GeForce
 633 RTX 2060—f-BRS-B took approximately 0.73s on
 634 average per seed points pair to segment each image
 635 (1296×864 image resolution) in the dataset, while
 636 BRS took 3.32s. This means that f-BRS-B is ap-

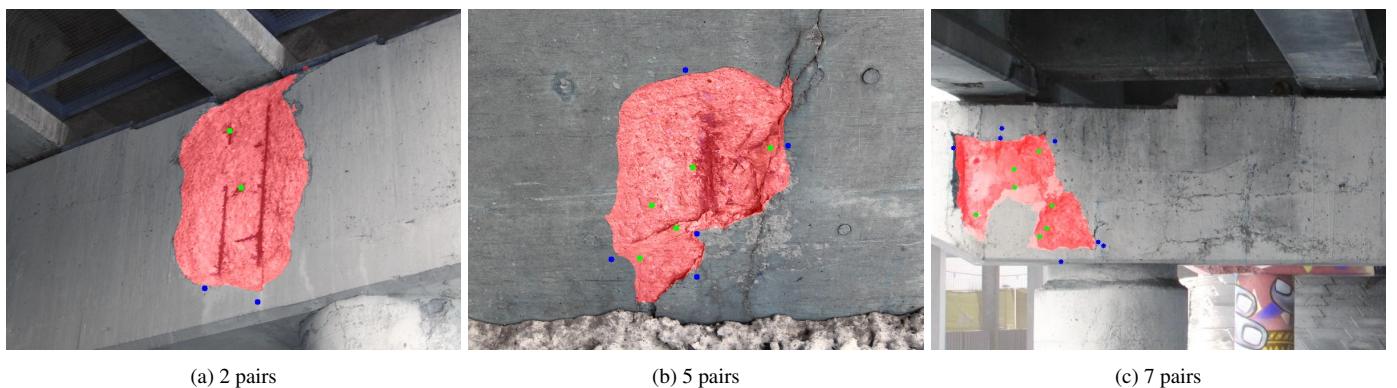


Figure 5: Sample f-BRS-B damage image segmentation results when area measurement accuracy is above 80%, for 2, 5, or 7 pairs of positive (inside spalling region) and negative (outside spalling region) seed points pairs. Note that these images are illustrations and were not captured from an XR device.

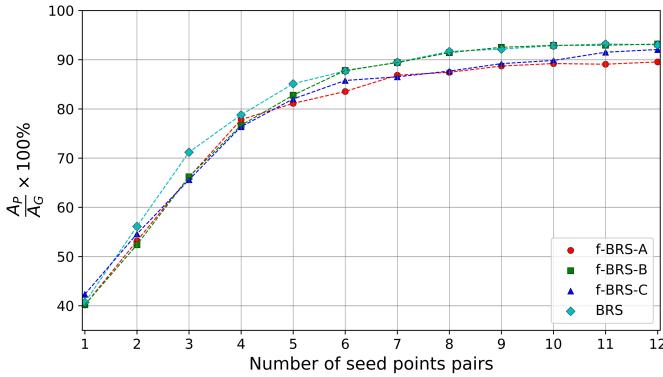


Figure 6: Evaluating area measurement accuracy as a function of the number of seed point pairs used for the segmentation models: f-BRS-A, f-BRS-B, f-BRS-C and BRS.

proximately 4.5 times faster than BRS while still maintaining a similar level of accuracy, which conformed to the findings of the original study [44]. This makes f-BRS-B the ideal implementation to choose for Xriv. Although f-BRS was developed to segment any visually distinct regions and its pretrained models were not trained on images having spalling-type damage, it can produce a reasonable performance for segmenting defects like spalling while maintaining a reasonable computation time. Note that since this experiment was performed using automated seed point selection, the accuracy of segmentation is lower than if we used human input to correct inaccurate segmentation results.

It should be noted that the number of seed points required depends on the unique situations, e.g., locations, shapes of damage regions, surface texture variation between damage and background, and ambiguous damage boundary (e.g., concrete spalling with flaking). This can be seen in Fig. 5. Thus, although we obtain over 80% accuracy using five pairs of seed points on average in this experiment, the ideal number of seed points required to obtain a satisfactory accuracy may vary depending on these conditions and the user's judgment.

4.2. HoloLens 2 Performance Evaluation

HL2, used as the XR device in this study, uses a SLAM algorithm to simultaneously build a spatial mesh of the user's environment while localizing the built-in camera's pose relative to the spatial mesh in real-time. The spatial computation component

of Xriv, outlined in Section 2.2, relies on the XR device's SLAM system as an input to compute the physical size of segmented damage regions. Therefore, an experiment was designed and conducted to evaluate the performance limits of HL2, which is directly related to the performance of Xriv.

Most of the SLAM algorithms utilize either pure visual feature matching and tracking, or a fusion of visual features and 3D depth/point cloud information. For a comprehensive survey of SLAM algorithms, readers are referred to [26, 27, 28, 29] and SLAM applications for infrastructure [15, 30, 16, 18, 31]. Due to the proprietary nature of HL2, details of the SLAM algorithm are not available in the public domain or in their documentation.

In general, all SLAM algorithms suffer to some degree from localization and mapping errors which accumulate over time due to uncertainties in the sensor measurements [54]. Localization error is the error associated with the position of the sensor within the global map and the mapping error is the error in the map features compared to the real scene. Due to the tightly coupled nature of SLAM, disambiguation of these two sources of errors can be challenging and these sources contribute to each other. Hence, both these errors affect the overall performance of Xriv and in lieu of such previous documentation regarding its accuracy, a separate experimental validation of HL2's performance was necessary.

The experiment was designed to evaluate the measurement error of HL2 as a function of the working distance (distance between the camera and measurement target). Different scene lighting conditions were also evaluated. Two 12 cm x 12 cm fiducial (ArUco) markers with a 15 cm center-to-center distance were attached on indoor (Fig. 7a) and outdoor (Fig. 7b) walls to investigate the effect of illumination [55]. The measurements are made possible using the built-in infrared-based depth sensor (Microsoft Azure Kinect [13]), whose performance is known to degrade under direct sunlight due to infrared interference [56].

The procedure was to capture images using HL2 at various distances from the wall. Here, the working distance is defined as the distance between the headset's built-in camera and the markers. The markers were positioned on the wall so that they match the



Figure 7: Evaluating the measurement accuracy of HL2 depending on a working distance under (a) indoor and (b) outdoor setting. Note that fiducial markers are used for automated detection.

same height as the headset’s camera. Experiments were conducted at working distances ranging from 0.5 m to 7.5 m (device specifications limit this to a maximum of 10 m [13]). An ArUco marker detector implemented in OpenCV [55] was applied to each image to detect the four corners of each marker in the images. The 3D world coordinates of the corner points were obtained by ray-casting to the spatial mesh, as described in Section 2.2, and the center of each marker was computed from its four corner locations. Finally, the distance between the center points was computed, and the measurement error was obtained by subtracting the true distance, 15 cm. The working distance was computed from the average distance between the camera center (obtained from its camera pose) and the centers of the two markers. This test was repeated both indoor and outdoor.

Results of 236 trials are shown in Fig. 8 as a scatter plot, where each point represents the distance measurement error at a certain working distance in either indoor or outdoor lighting conditions. The gray area indicates a 95% confidence interval (computed using a sliding window with a 2 m working distance).

In Fig. 8, the graph shows that the average measurement errors and their spread increase as the working distances increase. We can qualitatively observe that the confidence interval of the data at the same confidence level widens when the working distance exceeds 3 m. This indicates that the sample variance of the data is dependent on the working distance and that the data may be heteroskedastic. A Breusch–Pagan test [57] is performed to check the null hypothesis that variance is independent of working distance. As a result, the null hypothesis was re-

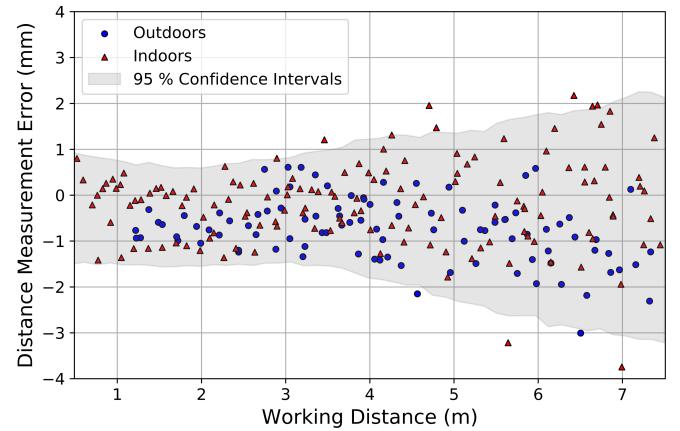


Figure 8: Scatter plot of distance measurement errors at different working distances between 0.5 m and 7.5 m under indoor and outdoor scene lighting conditions. Note that the gray area shows the 95% confidence interval (sliding window = 2 m).

jected with a p-value of 2.11×10^{-7} , which means that the variance increases as a function of working distance, which translates to the measurement accuracy dependence on working distance. This is not surprising because, as the working distances increase, the small-angle estimation errors in the camera matrix and/or the small offset of the spatial mesh can increase measurement errors, and the marker detection errors in pixels contribute to larger measurement errors.

Table 1 shows the mean error and spread at 95% confidence for both indoor and outdoor measurements, as well the aggregate of both cases. The mean

Table 1: Aggregate statistics of distance measurement error under indoor and outdoor settings (95% confidence interval).

Location	Mean error \pm spread
Indoor	0.16 ± 1.81 mm
Outdoor	0.76 ± 1.34 mm
Indoor & Outdoor	0.41 ± 1.73 mm

errors for indoor and outdoor cases are very small (below 1 mm), which indicates that HL2’s SLAM tracking system does not suffer from any significant drift regardless of brightness conditions and depth sensor interference from direct sunlight is not a significant factor for SLAM accuracy. The spread for both conditions is also relatively low (<2 mm), allowing consistent measurements. These experiments show that HL2 can produce sub-centimeter accuracy



(a)

(b)

Figure 9: Overview of test spalling regions at a bridge abutment (a) underneath bridge (SA) and (b) outside bridge (SB).

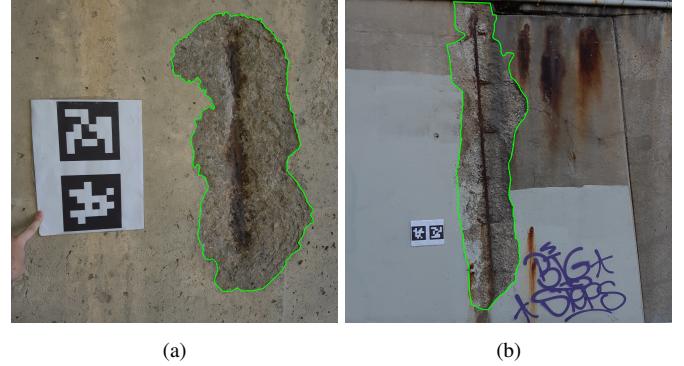
for both indoor and outdoor lighting conditions at working distances below 7.5 m. This level of accuracy is generally acceptable and sufficient for vision-based inspection applications [2, 3].

In this experiment, the difference in brightness between indoor and outdoor settings was compared qualitatively for the purpose of demonstrating that HL2 performs with reasonable accuracy under direct sunlight. Future studies can quantitatively measure the minimum and maximum amounts of brightness (measured in lumens) required to use HL2 with sufficient accuracy for common inspection scenarios where there may be either insufficient or excess amount of lighting.

4.3. Field Test

XRIV was tested at an in-service bridge in Southern Ontario, Canada. In this experiment, two spalling locations on the bridge were chosen from the abutment at inside (Fig. 9a) and outside (Fig. 9b) the testing bridge. Two spalling regions, labeled as SA and SB, have different sizes and approximately 0.5 m and 2 m in a vertical length, respectively. Both regions are present in locations that typically experience different lighting conditions. The working distance for each region was selected so that the entire region can be fit in a single image frame. SA and SB were captured from 1 m and 5 m away from the abutment, respectively.

The ground-truth areas of SA and SB were obtained first to be compared with the measurements



(a)

(b)

Figure 10: Measurement of ground-truth spalling areas for (a) SA and (b) SB.

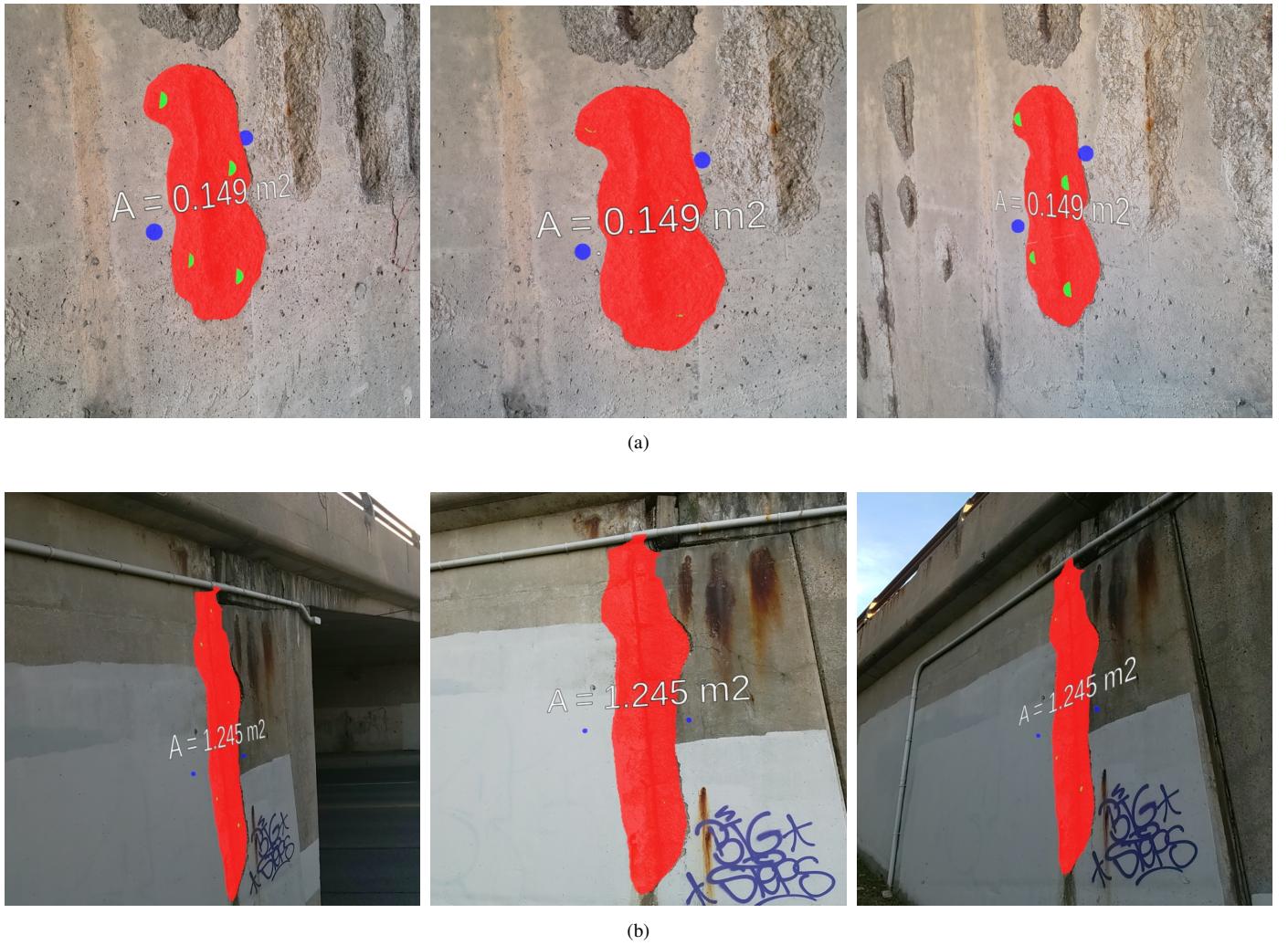
using XRIV. Since the spalling regions have irregular shapes and are often placed at inaccessible locations, like SB, an image-based approach was implemented to obtain their areas, instead of using other tools such as measuring tape. Fiducial markers with a known size were placed near the spalling region, and the entire view of the spalling region with the marker was captured as an image, shown in Fig. 10. Then, using a 4-point homography algorithm, a transformation matrix was obtained to transform all image pixel coordinates to corresponding physical coordinates [50]. The boundary of spalling on the images was manually segmented, which is shown as an outline in Fig. 10a and Fig. 10b for SA and SB, respectively. Then, the pixel coordinates of the boundary in the images were transformed to the physical scale coordinates, and then these coordinates were used



(a)

(b)

Figure 11: Seed points selected on SA in (a) and SB in (b): The points inside and outside spalling regions indicate positive and negative seed points, respectively. These images were captured on the HoloLens 2.



(a)

(b)

Figure 12: Visualization of the segmented regions for SA in (a) and SB in (b) and their areas, captured on the HoloLens 2. Since graphic objects are anchored to a spatial mesh of the scenes, the holographic objects are rendered based on the user's viewpoints.

819 to calculate the area of each spalling region using
820 the Shoelace formula explained in Section 2.3. The
821 ground-truth areas of SA and SB are presented in Ta-
822 ble 2.

823 The performance of XRIV was tested in both real-
824 time and offline modes. For the real-time mode,
825 seed points were selected and fed as inputs into the
826 XR interface for interactive segmentation of spalling
827 boundary (see Section 3). As discussed in Sec-
828 tion 4.1, the number of initial seed points required
829 to obtain a satisfactory result is variable and is based
830 on how difficult it is to distinguish the damage re-
831 gion from the background, so it is up to the user to
832 decide how many seed points are needed until a visu-
833 ally satisfactory result was obtained. Fig. 11 shows
834 the scenes of the XR interface after selecting seed
835 points from SA in Fig. 11a and SB in Fig. 11b. Then,
836 the images including a full view of the spalling re-
837 gions and the selected seed points were captured and
838 sent to a mobile computer (equipped with a GPU:
839 NVIDIA GeForce RTX 2060, which is the compu-
840 tational server outlined in Fig. 3) to segment the
841 spalling. In XRIV, the user can adjust the place-
842 ment and number of seed points through trial and
843 error to obtain the best segmentation results. How-
844 ever, in this experiment, accurate boundaries for SA
845 and SB could be obtained using six seed points (four
846 inside and two outside spalling regions) without re-
847 fining the process. Lastly, the areas of the segmented
848 regions for SA and SB were computed, and a holo-
849 graphic overlay of the segmented region with a text
850 overlay of the area was displayed to the XR interface,
851 as shown for both spalling regions in Fig. 12. Note
852 that the XR interface determines the shapes and posi-
853 tions of the graphic objects based on the user’s head-
854 set position and eye gaze relative to the spatial mesh.
855 Thus, the holograms including the color-filled region
856 and the text look as if they are stationary and are part
857 of the real scene.

858 In this experiment, HL2 was connected to the
859 computer using a Wi-Fi hotspot. It took around
860 2-3 seconds (or 0.5s per seed) to run the GPU-
861 accelerated f-BRS image segmentation on an im-
862 age with 908×511 resolution when six seed points
863 were provided. The speed of the segmentation in
864 XRIV is comparable to the one reported in the orig-
865 inal f-BRS study, which can be as fast as 0.32 sec-

866 onds per seed point when a GTX 1080 Ti GPU was
867 used [44]. Also, it took around 5-10 seconds for the
868 user to select six seed points through the XR inter-
869 face. The area computation runs in real-time, ex-
870 cluding the time taken to operate the holographic in-
871 terface, which is around 1 second. Thus, it takes ap-
872 proximately a total of 8-14 seconds to measure the
873 area of a single defect region in the real-time mode,
874 unless the user refines the segmented region which
875 may add a couple of seconds.

876 An offline mode was also used to estimate the
877 area. Images of SA and SB were captured and
878 stored on the device along with a spatial mesh of
879 the scene. Then, the boundaries of SA and SB in
880 the images were manually segmented, although au-
881 tomated segmentation algorithms could be used for
882 this step. Basically, the measurement errors from the
883 real-time and offline modes are identical if the seg-
884 mented boundaries of the damage regions are equal
885 in both modes. Thus, the measurement accuracy of
886 both modes was ascertained from this testing to eval-
887 uate the error resulting from the interactive segmen-
888 tation, assuming manual segmentation is the ground-
889 truth. Finally, the areas of SA and SB were computed
890 based on the procedure described in Section 2.3. The
891 results of the experiment are presented in Table 2,
892 which shows the damage size measurements for the
893 real-time and offline modes, compared to ground-
894 truth measurements.

Table 2: Damage size measurements for tested spalling re-
gions using XRIV in real-time and offline modes, compared to
ground-truth measurements.

Target damage	SA	SB
Real-time mode	0.149 m^2	1.245 m^2
Offline mode	0.143 m^2	1.176 m^2
Ground-truth	0.143 m^2	1.143 m^2
Error (real-time)	4.2%	8.9%
Error (offline)	<0.7%	2.9%

895 The results show that the real-time mode produces
896 larger errors compared to the offline mode, which is
897 expected due to the additional error caused by im-
898 age segmentation since the offline results were seg-
899 mented manually. Another result that can be ob-
900 served is that the error for SB is larger than SA,

which is likely because the working distance for SB is larger than SA, which are approximately 5 m and 1 m, respectively. These results are consistent with HL2's performance evaluation experiment in Section 4.2 that shows that distance measurement is dependent on working distance. This experiment also shows that Xriv generally produces errors that are within 10% of the ground-truth, for both real-time and offline modes, which makes it suitable for vision-based inspection applications [2, 3].

5. Conclusion

In this study, a new vision-based damage quantification system, Xriv, using XR technology is presented. Xriv allows the user to accurately measure the size of remote damage on structures by leveraging the built-in sensors and spatial mapping capabilities of a wearable XR device. The novelty of Xriv is that the user can interact with the damage detection process in real-time and obtain quantitative damage information. A DNN-based interactive segmentation algorithm, f-BRS, was deployed to segment damage regions in images using seed points provided by the user through the XR interface. The physical size of the segmented damage region is geometrically evaluated using a ray-casting algorithm to obtain 3D coordinates of the damage regions from pixel coordinates. The performance of Xriv is experimentally demonstrated by conducting a field experiment at an in-service bridge in the case where spalling is the target damage. This system has been shown to successfully estimate the size of spalling damage with less than 10% errors for two different test cases with latency times of 2-3s (or 0.5s per seed point) due to image processing. The Xriv system demonstrates how XR devices can allow expert users to utilize human-machine collaboration to achieve desired outcomes such as improved accuracy, robustness, real-time analysis, and immersive visualization of data. It is expected that Xriv will become an exemplary implementation of XR technology for human-assisted visual inspection.

Acknowledgments

We acknowledge the support from Rogers and the Ontario Centers of Excellence via *Voucher for Inno-*

vation and Productivity Program [OCE-34028] and Natural Sciences and Engineering Research Council of Canada (NSERC), [RGPIN-2020-03979].

References

- [1] C. McNally, B. Ferreira, D. L. A. Gordon, The Canadian infrastructure report card, Canada The State of the Federation 2015 (2015) 27–44.
- [2] MTO, Ontario Structure Inspection Manual (OSIM), Ministry of Transportation Ontario, 2008.
- [3] AASHTO, Manual for bridge element inspection, American Association of State Highway and Transportation Officials, 2019.
- [4] M. Moore, B. M. Phares, B. Graybeal, D. Rolander, G. Washer, J. Wiss, et al., Reliability of visual inspection for highway bridges, volume i, Tech. rep., Turner-Fairbank Highway Research Center (2001).
- [5] C. M. Yeum, S. J. Dyke, Vision-based automated crack detection for bridge inspection, Computer-Aided Civil and Infrastructure Engineering 30 (10) (2015) 759–770. doi:10.1111/mice.12141.
- [6] B. F. Spencer Jr, V. Hoskere, Y. Narazaki, Advances in computer vision-based civil infrastructure inspection and monitoring, Engineering 5 (2) (2019) 199–222. doi:10.1016/j.eng.2018.11.030.
- [7] E. McLaughlin, N. Charron, S. Narasimhan, Automated defect quantification in concrete bridges using robotics and deep learning, Journal of Computing in Civil Engineering 34 (5) (2020) 04020029. doi:10.1061/(ASCE)CP.1943-5487.0000915.
- [8] D. Li, Q. Xie, X. Gong, Z. Yu, J. Xu, Y. Sun, J. Wang, Automatic defect detection of metro tunnel surfaces using a vision-based inspection system, Advanced Engineering Informatics 47 (2021) 101206. doi:10.1016/j.aei.2020.101206.
- [9] P. Milgram, F. Kishino, A taxonomy of mixed reality visual displays, IEICE TRANSACTIONS on Information and Systems 77 (12) (1994) 1321–1329.
- [10] R. K. Ganesan, Mediating human-robot collaboration through mixed reality cues, Ph.D. thesis, Ph. D. Dissertation. Arizona State University (2017).
- [11] M. Khatib, K. Al Khudir, A. De Luca, Human-robot contactless collaboration with mixed reality interface, Robotics and Computer-Integrated Manufacturing 67 (2021) 102030. doi:10.1016/j.rcim.2020.102030.
- [12] Microsoft, HoloLens (1st gen) hardware (2019).
- [13] Microsoft, HoloLens 2 hardware (2020).
- [14] Magic Leap, Magic Leap 1 (2019).
- [15] A. Kim, R. M. Eustice, Real-time visual slam for autonomous underwater hull inspection using visual saliency, IEEE Transactions on Robotics 29 (3) (2013) 719–733. doi:10.1109/TRO.2012.2235699.
- [16] H. Peel, S. Luo, A. Cohn, R. Fuentes, Localisation of a mobile robot for bridge bearing inspec-

- tion, *Automation in Construction* 94 (2018) 244–256. doi:10.1016/j.autcon.2018.07.003.
- [17] V. Hoskere, J.-W. Park, H. Yoon, B. F. Spencer Jr, Vision-based modal survey of civil infrastructure using unmanned aerial vehicles, *Journal of Structural Engineering* 145 (7) (2019) 04019062. doi:10.1061/(ASCE)ST.1943-541X.0002321.
- [18] N. Charron, E. McLaughlin, S. Phillips, K. Goorts, S. Narasimhan, S. L. Waslander, Automated bridge inspection using mobile ground robotics, *Journal of Structural Engineering* 145 (11) (2019) 04019137. doi:10.1061/(ASCE)ST.1943-541X.0002404.
- [19] S. Chen, D. F. Laefer, E. Mangina, S. I. Zolanivari, J. Byrne, Uav bridge inspection through evaluated 3d reconstructions, *Journal of Bridge Engineering* 24 (4) (2019) 05019001. doi:10.1061/(ASCE)BE.1943-5592.0001343.
- [20] Y.-z. Lin, Z.-h. Nie, H.-w. Ma, Structural damage detection with automatic feature-extraction through deep learning, *Computer-Aided Civil and Infrastructure Engineering* 32 (12) (2017) 1025–1046. doi:10.1111/mice.12313.
- [21] Y.-J. Cha, W. Choi, G. Suh, S. Mahmoudkhani, O. Büyüköztürk, Autonomous structural visual inspection using region-based deep learning for detecting multiple damage types, *Computer-Aided Civil and Infrastructure Engineering* 33 (9) (2018) 731–747. doi:10.1111/mice.12334.
- [22] C. M. Yeum, J. Choi, S. J. Dyke, Automated region-of-interest localization and classification for vision-based visual assessment of civil infrastructure, *Structural Health Monitoring* 18 (3) (2019) 675–689. doi:10.1177/1475921718765419.
- [23] S. Li, X. Zhao, G. Zhou, Automatic pixel-level multiple damage detection of concrete structure using fully convolutional network, *Computer-Aided Civil and Infrastructure Engineering* 34 (7) (2019) 616–634. doi:10.1111/mice.12433.
- [24] B. Kim, S. Cho, Automated multiple concrete damage detection using instance segmentation deep learning model, *Applied Sciences* 10 (22) (2020) 8008. doi:10.3390/app10228008.
- [25] M. Anitha, R. Hemalatha, S. Radha, A survey on crack detection algorithms for concrete structures, in: *Advances in Smart System Technologies*, Springer, 2021, pp. 639–654. doi:10.1007/978-981-15-5029-4_53.
- [26] T. Taketomi, H. Uchiyama, S. Ikeda, Visual slam algorithms: a survey from 2010 to 2016, *IPSJ Transactions on Computer Vision and Applications* 9 (1) (2017) 16. doi:10.1186/s41074-017-0027-2.
- [27] M. Filipenko, I. Afanasyev, Comparison of various slam systems for mobile robot in an indoor environment, in: *2018 International Conference on Intelligent Systems (IS)*, IEEE, 2018, pp. 400–407. doi:10.1109/IS.2018.8710464.
- [28] B. Huang, J. Zhao, J. Liu, A survey of simultaneous localization and mapping, *arXiv preprint arXiv:1909.05214* (2019).
- [29] C. Debeunne, D. Vivet, A review of visual-lidar fusion based simultaneous localization and mapping, *Sensors* 20 (7) (2020) 2068. doi:10.3390/s20072068.
- [30] Z. Shang, Z. Shen, Real-time 3d reconstruction on construction site using visual slam and uav, *arXiv preprint arXiv:1712.07122* (2017). doi:10.1061/9780784481264.030.
- [31] S. Jung, D. Choi, S. Song, H. Myung, Bridge inspection using unmanned aerial vehicle based on hg-slam: Hierarchical graph-based slam, *Remote Sensing* 12 (18) (2020) 3022. doi:10.3390/rs12183022.
- [32] X. Chen, J. Li, S. Huang, H. Cui, P. Liu, Q. Sun, An automatic concrete crack-detection method fusing point clouds and images based on improved otsu's algorithm, *Sensors* 21 (5) (2021) 1581. doi:10.3390/s21051581.
- [33] A. Webster, S. Feiner, B. MacIntyre, W. Massie, T. Krueger, Augmented reality in architectural construction, inspection and renovation, in: *Proc. ASCE Third Congress on Computing in Civil Engineering*, Vol. 1, 1996, p. 996.
- [34] R. T. Azuma, A survey of augmented reality, *Presence: Teleoperators & Virtual Environments* 6 (4) (1997) 355–385. doi:10.1162/pres.1997.6.4.355.
- [35] J. Carmigniani, B. Furht, M. Anisetti, P. Ceravolo, E. Damiani, M. Ivkovic, Augmented reality technologies, systems and applications, *Multimedia tools and applications* 51 (1) (2011) 341–377. doi:10.1007/s11042-010-0660-6.
- [36] A. H. Behzadan, S. Dong, V. R. Kamat, Augmented reality visualization: A review of civil infrastructure system applications, *Advanced Engineering Informatics* 29 (2) (2015) 252–267. doi:10.1016/j.aei.2015.03.005.
- [37] V. R. Kamat, S. El-Tawil, Evaluation of augmented reality for rapid assessment of earthquake-induced building damage, *Journal of computing in civil engineering* 21 (5) (2007) 303–310. doi:10.1061/(ASCE)0887-3801(2007)21:5(303).
- [38] R. Napolitano, Z. Liu, C. Sun, B. Glisic, Combination of image-based documentation and augmented reality for structural health monitoring and building pathology, *Frontiers in Built Environment* 5 (2019) 50. doi:10.3389/fbui.2019.00050.
- [39] F. Moreu, B. Bleck, S. Vemuganti, D. Rogers, D. Mascarenas, Augmented reality tools for enhanced structural inspection, *Structural Health Monitoring* 2 (2017). doi:10.12783/shm2017/14221.
- [40] D. Maharjan, M. Agüero, D. Mascarenas, R. Fierro, F. Moreu, Enabling human-infrastructure interfaces for inspection using augmented reality, *Structural Health Monitoring* (2020). doi:10.1177/1475921720977017.
- [41] E. Karaaslan, U. Bagci, F. N. Catbas, Artificial intelligence assisted infrastructure assessment using mixed reality systems, *Transportation Research Record* 2673 (12) (2019) 413–424. doi:10.1177/0361198119839988.
- [42] Microsoft, MixedReality-SpectatorView (2019).

- 1110 [43] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the
 1111 IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.
 1112
- 1113 [44] K. Sofiuk, I. Petrov, O. Barinova, A. Konushin, f-brs: Re-thinking backpropagating refinement for interactive seg-
 1114 mentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp.
 1115 8623–8632.
 1116
- 1117 [45] J. E. van Engelen, H. H. Hoos, A survey on semi-
 1118 supervised learning, *Machine Learning* 109 (2) (2020)
 1119 373–440. doi:10.1007/s10994-019-05855-6.
 1120
- 1121 [46] N. Xu, B. Price, S. Cohen, J. Yang, T. S. Huang, Deep
 1122 interactive object selection, in: Proceedings of the IEEE
 1123 Conference on Computer Vision and Pattern Recognition,
 1124 2016, pp. 373–381.
 1125
- 1126 [47] W.-D. Jang, C.-S. Kim, Interactive image segmenta-
 1127 tion via backpropagating refinement scheme, in: Pro-
 1128 ceedings of the IEEE Conference on Computer Vi-
 1129 sion and Pattern Recognition, 2019, pp. 5297–5306.
 1130 doi:10.1109/CVPR.2019.00544.
 1131
- 1132 [48] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam,
 1133 Encoder-decoder with atrous separable convolution for
 1134 semantic image segmentation, in: Proceedings of the Eu-
 1135 ropean conference on computer vision (ECCV), 2018, pp.
 1136 801–818. doi:10.1007/978-3-030-01234-2_49.
 1137
- 1138 [49] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, J. Malik,
 1139 Semantic contours from inverse detectors, in: 2011 In-
 1140 ternational Conference on Computer Vision, IEEE, 2011,
 1141 pp. 991–998. doi:10.1109/ICCV.2011.6126343.
 1142
- 1143 [50] R. Hartley, A. Zisserman, Multiple view geometry in
 1144 computer vision, Cambridge university press, 2003.
 1145 doi:10.1017/CBO9780511811685.
 1146
- 1147 [51] C. Ericson, Real-time collision detection, CRC Press,
 1148 2004. doi:10.1201/b14581.
 1149
- 1150 [52] Y. Lee, W. Lim, Shoelace formula: Connecting the
 1151 area of a polygon and the vector cross product,
 1152 The Mathematics Teacher 110 (8) (2017) 631–636.
 1153 doi:10.5951/mathteacher.110.8.0631.
 1154
- 1155 [53] S. Suzuki, et al., Topological structural analysis of dig-
 1156 itized binary images by border following, *Computer vi-
 1157 sion, graphics, and image processing* 30 (1) (1985) 32–46.
 1158 doi:10.1016/0734-189X(85)90016-7.
 1159
- 1160 [54] R. Mur-Artal, J. D. Tardós, Orb-slam2: An open-source
 1161 slam system for monocular, stereo, and rgb-d cameras,
 1162 *IEEE Transactions on Robotics* 33 (5) (2017) 1255–1262.
 1163 doi:10.1109/TRO.2017.2705103.
 1164
- 1165 [55] Itseez, Open source computer vision library, <https://github.com/itseez/opencv> (2020).
 1166
- 1167 [56] M. R. Andersen, T. Jensen, P. Lisouski, A. K. Mortensen,
 1168 M. K. Hansen, T. Gregersen, P. Ahrendt, Kinect depth
 1169 sensor evaluation for computer vision applications, Tech.
 1170 rep., Aarhus University (2012).
 1171
- 1172 [57] T. S. Breusch, A. R. Pagan, A simple test for het-
 1173 eroscedasticity and random coefficient variation, *Econo-
 1174 metrica: Journal of the Econometric Society* (1979)

1287–1294doi:10.2307/1911963.

1166