

Robust Enhancement of Depth Images from Depth Sensors

ABM Tariqul Islam^a, Christian Scheel^a, Renato Pajarola^b, Oliver Staadt^{a,*},

^aVisual Computing Lab, University of Rostock, Germany.
^bVisualization and MultiMedia Lab, University of Zürich, Switzerland.

Abstract

In recent years, depth cameras (such as Microsoft Kinect and ToF cameras) have gained much popularity in computer graphics, visual computing and virtual reality communities due to their low price and easy availability. While depth cameras (e.g. Microsoft Kinect) provide RGB images along with real-time depth information at high frame rate, the depth images often suffer from several artifacts due to inaccurate depth measurement. These artifacts highly degrade the visual quality of the depth frames. Most of these artifacts originate from two main sources – the missing/invalid depth values and fluctuating valid depth values on the generated contents. In this paper, we propose a new depth image enhancement method, for the contents of depth cameras, which addresses these two main sources of artifacts. We introduce a robust 1D Least Median of Squares (1D LMedS) approach to estimate the depth values of those pixels which have missing/invalid depth values. We use a sequence of frames to look for invalid depth values (considered as outliers), and finally, replace those values with stable and more plausible depth values. By doing so, our approach improves the unstable nature of valid depth values in captured scenes that is perceived as flickering. We use self-recorded and reference datasets along with reference methods to evaluate the performance of our proposed 1D LMedS. Experimental results show improvements both for static and moving parts of a scene.

Keywords: Robust filtering, RGB-D sensor, Depth image, 3D imaging, Depth image enhancement, Virtual reality.

1. Introduction

Nowadays, researchers often use depth cameras (e.g. Microsoft Kinect, ToF cameras) to generate 3D contents. Due to their low cost and easy availability, several areas in computer graphics (e.g. visual computing, human-computer interaction, virtual reality, remote collaboration) adopted these cameras as the primary acquisition device. The sensors of these cameras capture a scene with acceptable resolution and speed, but the captured depth data exhibit considerable amount of artifacts (e.g., flickering artifact and random black spots as shown in Figure 1). Most of these artifacts originate from two primary sources—invalid depth values and fluctuating valid depth values on the generated contents. Due to these two prime sources of artifacts, we often see highly degraded 3D contents whose quality is often too low for many applications such as telepresence,

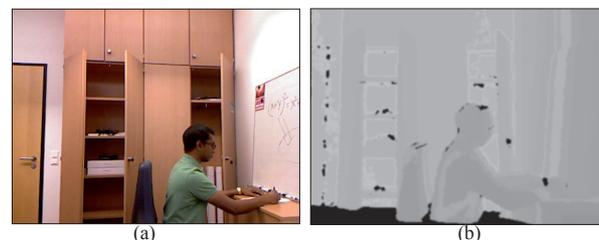


Figure 1: RGB-D frame from a Kinect sensor – (a) color image, (b) depth image. The artifacts (e.g. the black spots) in (b) occurs due to the invalid depth values reported by Kinect RGB-D sensor.

object tracking in a scene, forensic analysis of crime scenes etc. In case of telepresence, a low quality 3D scene capturing limits the sensation of a natural presence of remote users on the local site [1]. Similarly, for tracking applications, the noise on an object's surface greatly deteriorates the tracking performance. Likewise, in forensic analysis of a crime scene, where every detail is crucial, a degraded scene might hamper the recovery of correct information [2].

Accordingly, researchers conducted a number of

*Corresponding author

Email addresses: tariqul.islam@uni-rostock.de (ABM Tariqul Islam), CScheel@gmx.de (Christian Scheel), pajarola@ifi.uzh.ch (Renato Pajarola), oliver.staadt@uni-rostock.de (Oliver Staadt)

works to enhance the depth frames from the depth cameras by addressing the aforementioned sources of artifacts. In those works, they propose to use different filters in combination with motion compensation and color image data (see e.g., [1, 3–17]). In this paper, we present a robust method of enhancing depth frames which is a continuation and extension of our previous works [18, 19]. Here we have extended our previous works by including a more detailed technical description of our proposed method, several additional datasets and reference works to evaluate our work and, quantitative and qualitative comparison analysis on real-world and artificial depth data. Moreover, we clarify the concept of general Least Median of Squares and the concept of how our method improves each pixels of RGB-D frames. Furthermore, we demonstrate how our method can be implemented to improve the 3D representation quality of a telepresence communication.

Our approach. Contrary to prior works, we propose the use of Least Median of Squares (LMedS)[20] and extend it to enhance both static and dynamic scenes. By adapting LMedS to be used as a 1D (one dimensional) LMedS, we go through a sequence of frames, detect the invalid and unstable depth values for each pixel, identify the outliers among them and finally, replace them with stable valid depth values. Our proposed 1D LMedS is computationally less expensive than higher dimensional LMedS and the application of 1D LMedS on the depth frames is efficient. Experimental results show improvements in filling the invalid depth values and stabilizing the valid values for static and dynamic parts of a scene.

Contribution. Our contribution towards enhancing the depth frames of depth cameras is as follows:

- A computationally inexpensive and efficient way of enhancing the depth images for real-time applications.
- Introducing 1D LMedS which fits right into the problem domain of depth image enhancement.
- Formulate the estimator, standard deviation and weight function from original LMedS to fit into 1D LMedS.
- Develop a fast data traversal procedure for 1D LMedS.
- Correctly identify the invalid depth values and fluctuating valid depth values, and stabilize flickering artifacts for both the static and dynamic scenes.
- Evaluate our method with local and reference datasets and with reference methods.

2. Related work

Over the past few years, we have seen a number of works on enhancing the depth images captured by depth cameras. Most of these works use filters, like median filters [1, 4], Kalman filters [5], guided image filters [6, 7, 21, 22], bilateral filters [8, 10–15], and sometimes a combination of median and bilateral filters [12, 23]. Beside the type of applied filter, the approaches can basically be categorized by the application of motion compensation, the covered depth sensor type, the inclusion of color image data, and realtime support, for example, by hardware acceleration. Most of these approaches utilize spatial or range filter kernels, while only very few of them consider the temporal noise component.

For instance, Maimone and Fuchs use a GPU-adapted median filter for enhancing the depth images [1]. Their method successfully fills the black spots on the surface of the objects, but it occasionally reports wrong interpolated values. Moreover, it also exhibits temporally inconsistent depth frames due to not considering temporal information. Matyunin et al. use motion information along with a median filter which generates good results for some image areas, but suffers mainly from noisy object boundaries [4]. Besides, it suffers from large computational overhead. Lai et al. propose to use a recursive median filter on the depth images from Kinect sensors. However, they mainly focus on improving object recognition accuracy rather than depth map accuracy [24].

Amamra and Aouf use a GPU-adaptive Kalman filter to enhance the depth images in real-time [5]. However, this method suffers when object’s surface is sharply reflected. Wasza et al. propose another GPU-based approach where a combination of normalized convolution and guided filter is used to improve the depth map in real-time, but it suffers from blurry object borders [25]. Some other works (e.g. [10, 11]) use bilateral filters for depth image enhancement. Chen et al. use a joint bilateral filter where color image is used as guidance image [11]. Their method enhances the depth images except for the image areas without color information. Yang et al. also enhance depth images by using bilateral filters [12]; but it is not suitable for dynamic scenes because it does not consider temporal information.

Hui and Ngan use a combination of 2D median and bilateral filters where the bilateral filter removes the artifacts introduced by the median filter [23]. Maimone et al. also use a combination of median and bilateral filters to enhance the depth images used in telepresence [26]. Although, their method enhances the surface of objects, it exhibits temporally inconsistent depth frames.

Telea uses an image inpainting method [27] which

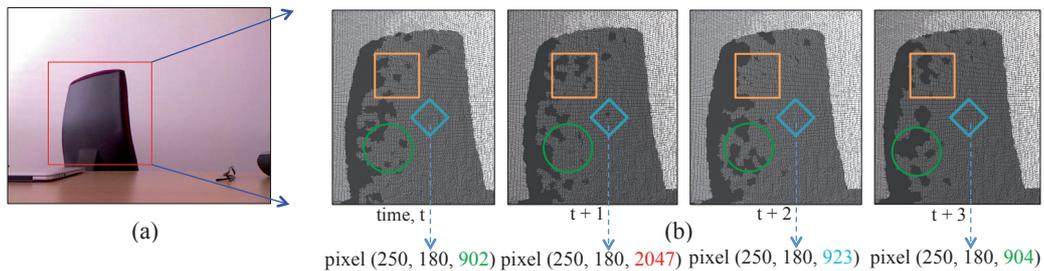


Figure 2: **RGB-D frames from Kinect with static objects** – (a) color image, (b) four consecutive depth frames of square marked area of (a). The depth values of one pixel inside the diamond marked area in (b) are depicted below the frames. We observe three types of depth value distortions: 2047 is an invalid value (reported occasionally), 923 is an outlier and the others are regular depth values fluctuating because of Gaussian noise.

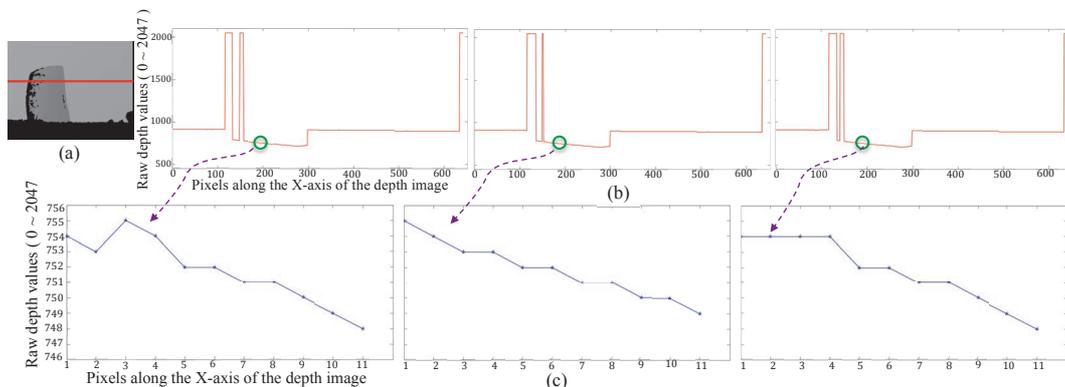


Figure 3: **Invalid and unstable valid depth values from Kinect sensor on a sequence of frames** – (a) depth image with the red marked scan line, (b) depth values of the scan line for 3 consecutive frames of (a), (c) depth values of the green circled marked area on the plots of (b).

127 considers both the depth and color information to fill
 128 in the missing spots [14]. However, it considers only
 129 the spatial relationship among the pixels, resulting in
 130 wrong interpolated values and noisy object boundaries.
 131 Camplani et al. use an iterative joint-bilateral filter
 132 which combines depth and color data by analyzing
 133 an edge-uncertainty map and pre-detected foreground
 134 regions [8, 28]. It performs well for static scenes,
 135 but is not suitable for dynamic scenes. Other works
 136 (e.g. [17, 29]) also focus on depth map processing, but
 137 these are based on stereo [29] or are specific to time-of-
 138 flight (ToF) sensors [17].

139 There are also other works which utilize color image
 140 data to enhance the quality of the corresponding depth
 141 images (e.g. [9, 30–32]). While most of these meth-
 142 ods use the Markov Random Field model or the auto re-
 143 gression model to enhance the depth images, they differ
 144 mainly in the design of regularization terms in the objec-
 145 tive function. The resulting depth images are quite good
 146 and contain almost no artifacts. Although, the energy
 147 minimization approaches, which these works use, are
 148 computationally expensive [30], recent computational

149 development would enable these methods to be used
 150 for realtime applications. There are also some works
 151 (e.g. [33]) which use color-depth image pairs and inter-
 152 polation based approach to refine and enhance the cor-
 153 responding depth images. The authors in [33] propose
 154 a combination of de-noise and interpolation method
 155 which performs quite well for stationary objects of a
 156 scene; it nicely refines the areas with artifacts in depth
 157 images.

158 We have also seen works (e.g. [17, 23, 34, 35]) which
 159 use the flow of information in consecutive color images.
 160 In [17], Kim et al. use a combination of spatial and tem-
 161 poral depth enhancement process that considers motion
 162 flow between consecutive color images to infer infor-
 163 mation about object motion in the corresponding depth
 164 images. However, they basically ignore this data in the
 165 dynamic parts of the depth images as they use it only to
 166 detect stationary parts. Based on this, they apply a bi-
 167 lateral filter to improve the quality which works well for
 168 static parts. Hui et. al. enhance depth images captured
 169 from a moving RGB-D system [23]. They also esti-
 170 mate the optical flow of consecutive color images. How-

171 ever, instead of building a temporal filter on top of the
172 obtained data, their method estimates additional depth
173 cues from the flow which are then combined with the
174 original depth images. Their method is intended for mo-
175 bile setups and not applicable to stationary cameras as
176 considered in our case. Another work [35] by Avetisyan
177 et. al. also use optical flow between successive frames
178 and they are able to get decent quality depth frames. Al-
179 though their method works in real-time, it suffers when
180 missing data or noise stays persistently in one region of
181 the depth image sequence. Besides, their method occa-
182 sionally yields invalid motion vectors which deteriorate
183 the estimated history of depth values; as a result, ghost-
184 ing artifacts are shown for fast moving objects.
185 Even though each of the mentioned works is designed
186 to solve different issues related to the depth data from
187 depth cameras and they perform well in their respective
188 applications areas, most of them lack in considering the
189 temporal inconsistency and often, they suffer from large
190 computational overhead. In this paper, we address the
191 temporal inconsistency issue of depth cameras.

192 3. Problem description

193 Depth cameras (such as Kinects and ToF cameras) ex-
194 hibit artifacts in the form of missing data, invalid depth
195 values and valid fluctuating depth values. Here we use
196 Kinect data to illustrate the artifacts, because major part
197 of the depth image enhancement works utilize Kinects
198 as acquisition device. Figure 1 shows the artifacts found
199 on Kinect’s depth frames. These artifacts occur because
200 the structured light sensor of a Kinect produces faulty
201 depth measurements. A Kinect captures both the color
202 and depth data of a scene. The depth of an object is
203 measured by projecting light patterns from an infrared
204 light source to the object space [36]. An infrared cam-
205 era receives the reflected light pattern from the surface
206 of the object and compares it against a reference pat-
207 tern. An estimated disparity image of the objects inside
208 a scene and corresponding depth calculations are ob-
209 tained by the differences between the captured patterns
210 and the reference pattern [8, 37]. The faulty depth mea-
211 surements of a Kinect might originate due to multiple
212 reflections, strong lighting conditions, scattering object
213 surface etc. [8, 36, 37].

214 Moreover, the depth measurements are also affected
215 by the instability over time, which results in flickering
216 spots [8]. Flickering occurs even on a scene with static
217 objects (see Figure 2); here, we can observe the unstable
218 nature of valid depth values for certain pixels. More-
219 over, Figure 3 (b) shows that a Kinect reports both in-
220 valid values (the spikes on the plots) and valid values

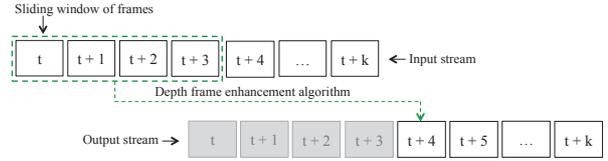


Figure 4: **Illustration of sliding window of frames in 1D LMedS** – number of frames inside the slide window is 4 and total number of frames is k .

221 for the same pixels in a sequence of frames, although
222 the objects are static. Besides, Figure 3 (c) shows that
223 the valid depth values, inside the circled marked areas
224 of Figure 3 (b), also differ for the same pixels in the
225 sequence of frames. We take this unstable nature into
226 consideration and apply 1D LMedS to replace the un-
227 stable and invalid values with valid stable depth values.

228 4. Least Median of Squares

The Least Median of Squares (LMedS) [20] is a statistical technique for robust regression of a p -dimensional sample set (x_i, y_i) , which minimizes the median of the squared residuals r_i^2 over a set of estimates θ in \mathbb{R}^p . Here, r_i and θ are expressed as in Equation 1 and Equation 2:

$$r_i(n) = y_i - (x_{i1}\hat{\theta}_1 + \dots + x_{i(p-1)}\hat{\theta}_{p-1} + x_{ip}\hat{\theta}_p) \quad (1)$$

$$\theta = (\theta_1, \theta_2, \dots, \theta_p)^t \quad (2)$$

229 Basically, a vector n represents a hyperplane that, in a
230 certain way, describes the linear dependency of y_i from
231 the $p-1$ variables in x_i . LMedS is known to be robust to
232 false matches and outliers while fitting equations to the
233 observed dataset. It estimates the data points by solv-
234 ing the nonlinear minimization problem which is $r_i(n)$.
235 The outliers from the dataset can be filtered out by us-
236 ing a robust standard deviation $\hat{\sigma}$ which uses the min-
237 imum median of $r_i(n)$. But, the LMedS estimator be-
238 comes computationally expensive when the dimension
239 p of sample set (x_i, y_i) gets bigger. In our approach,
240 we propose to use a one dimensional LMedS which is a
241 good fit for the nature of our problem domain and its im-
242 plementation in such problem domain is also efficient.

243 5. Proposed strategy

244 To enhance the depth frames by filling the invalid
245 depth values $d_{in.vld}$ and unstable valid depth values $d_{u.vld}$
246 with stable valid depth values $d_{s.vld}$, we introduce a 1D
247 LMedS approach which is robust to outliers. We apply

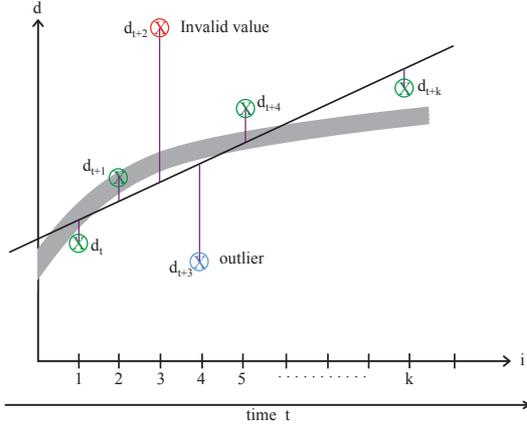


Figure 5: **Illustration of 1D LMedS on one depth pixel in a sliding window.** The thick gray line shows the real distance of a scene object over time. The checked circles represent depth values d_{t+i-1} obtained in k consecutive frames. The thin black line approximates the gray line by calculating the estimator m i.e. by minimizing the median of the absolute distances to the depth values, visualized by vertical purple lines. The red marked *invalid value* at $i = 3$ and the turquoise marked *outlier* at $i = 4$ are detectable by our 1D LMedS.

248 the 1D LMedS on a per pixel basis over a sequence of
 249 frames inside a temporal sliding window. Let us con-
 250 sider, for example, that we have k consecutive depth
 251 frames and we take n frames in a sliding window (see
 252 Figure 4). Due to Kinect’s nature of depth value gener-
 253 ation, we get different depth values d_i for any single pixel
 254 for the n previous consecutive frames inside the sliding
 255 window. We apply the 1D LMedS on this set of depth
 256 values which vary along the temporal domain.

257 The name 1D LMedS (one dimensional LMedS)
 258 comes from the concept that the depth values d_i of a
 259 certain pixel (x,y) changes only in one dimension that
 260 is along the temporal domain i.e. over the consecutive
 261 frames. Our goal is to locate the invalid and unstable
 262 depth values (which we consider as outliers) from the
 263 set of depth values of each pixel and replace them with
 264 a stable valid depth value. An illustration of 1D LMedS
 265 on k consecutive frames is depicted in Figure 5. Our
 266 depth enhancement strategy is based on three principle
 267 steps. For each pixel inside the sliding window (i.e. for
 268 the n consecutive frames):

- 269 • detect if there are invalid depth values $d_{in.vld}$.
- 270 • detect if the valid depth values d_{vld} are fluctuating
 271 from one frame to another.
- 272 • if the d_{vld} are fluctuating, find the outliers among
 273 the d_{vld} , remove them and process the inliers to get
 274 a final stable depth value $d_{s.vld}$ for that pixel. Also
 275 replace the $d_{in.vld}$ with the $d_{s.vld}$.

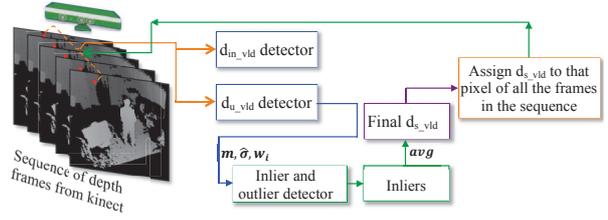


Figure 6: **Block diagram of the processing steps of 1D LMedS** – the red dots in the frame sequence represent one pixel of the frames.

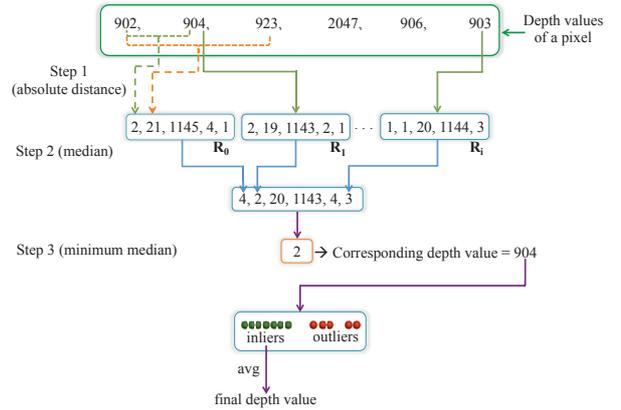


Figure 7: **Illustration of 1D LMedS for one pixel in the n consecutive frames** – the values on the top green box are the raw depth values obtained for one pixel over the frames inside the sliding window.

276 The processing steps of 1D LMedS are depicted in
 277 Figure 6. For detecting $d_{in.vld}$, we look for the pixels
 278 which have invalid depth values; the depth value of pix-
 279 els with invalid depth can be found in [8, 37]. To detect
 280 the unstable depth values $d_{u.vld}$ from the set of valid
 281 depth values d_{vld} for the n frames, we use the standard
 282 deviation $\hat{\sigma}$ of the depth values for the n frames and
 283 then compare with certain constraints (stated in Equa-
 284 tion 4) to remove them as outliers. A detailed descrip-
 285 tion of these processes is explained in the following
 286 paragraphs.

287 To remove the outliers and fill the invalid depth values
 288 for a certain pixel, we calculate the estimator m which
 289 is the minimum of the medians of absolute residuals R_i
 290 for each depth value d_i of the set of depth values for n
 291 consecutive frames. In 1D LMedS, we use the absolute
 292 difference for getting the residuals R_i among the depth
 293 values d_i ; in [20], R_i^2 is used to avoid negative values.

294 The whole calculation process is illustrated in Fig-
 295 ure 7. Here, in step 1, for a set of depth values d_i , we
 296 calculate absolute difference for each depth value to the
 297 rest of the depth values for that pixel. So, we obtain a set
 298 of residuals for each depth values of d_i . In Figure 7, the

299 first set of absolute distances i.e. $R_0 = \{2, 21, 1145, 4, 1\}$
 300 are calculated for the first depth value 902, R_1 for the
 301 second depth value 904 and so on. In step 2, we calcu-
 302 late the medians from each set of residuals R_i and then,
 303 in step 3, we calculate the 1D LMedS estimator m which
 304 is the minimum of the medians. Equation 3 shows the
 305 formula to calculate m .

$$R_i(d_i) = \{ |d_i - d_j|, \quad \text{if } i \neq j \} \\ (i,j:1,\dots,n) \quad (3)$$

$$m = \min_{(i:1,\dots,n)} \text{median } R_i(d_i)$$

After that, we retrieve the corresponding depth value whose index matches the index of m . Here, in step 3, we can see that the value of m is 2; which originates from the second set of residuals R_1 . As R_1 originates from the second depth value 904, we use the value of m and the retrieved depth value (904) for the rest of the process of 1D LMedS. We calculate the standard deviation $\hat{\sigma}$ and weight w_i , which decide if a certain depth value is an inlier or outlier, according to the formula stated in Equation 4. After we obtain inliers and outliers, we finally obtain the stable valid depth value by taking the average of the inliers.

$$\hat{\sigma} = 1.4826 \left[1 + \frac{5}{n-1} \right] m \quad (4)$$

$$w_i = \begin{cases} 1, & \text{if } r \leq 2.5\hat{\sigma} \\ 0, & \text{otherwise.} \end{cases}$$

306 The constant 1.4826 of Equation 4 is a coefficient to
 307 achieve the optimal efficiency in the presence of Gaus-
 308 sian noise, $5/(n-1)$ is used to compensate the effect of
 309 using a small set of data and the constant 2.5 is used with
 310 $\hat{\sigma}$ because in a Gaussian situation there will be very few
 311 residuals larger than $2.5\hat{\sigma}$; for a detailed explanation of
 312 the coefficients, please look at pages 17, 29, 44, 131,
 313 132 and 202 in [38]. However, this process needs quite
 314 some time to traverse for all the depth values of each
 315 pixel for these n frames. Because, in Figure 7, we can
 316 see that to obtain the absolute distances R_i s, we need to
 317 check each depth value to the rest of the depth values to
 318 get the R_i s which takes quite some time. Therefore, to
 319 reduce the calculation of R_i s, we adapt a fast data traver-
 320 sal procedure to fit 1D LMedS.

321 **Fast data traversal procedure for 1D LMedS.** The
 322 fast data traversal procedure is stated in Algorithm 1.
 323 We apply this procedure for each of the depth values d_i
 324 of every pixel of the n frames. Algorithm 1 shows the
 325 pseudocode for fast traversal for one pixel over the n
 326 frames. Here, the fast traversal occurs due to the faster

Algorithm 1 Fast traversal for the depth values of one pixel

```

1: procedure
2:   Define MAX 10000
3:   int idx = 0, n > n = no. of frames in the window
4:   Float diff, depthVal[n], preFinalDepthVal
   > depthVal[n] contains the unsorted depth values
   of a pixel
5:   depthVal_S[n]=sorted depth values of depthVal[n]
   > values of depthVal_S[n] are in ascending order
6:   min ← Float MAX
7:   for int i = 0 to  $\frac{n-1}{2}$  do
8:     diff = depthVal_S[ $i + \frac{n}{2}$ ] - depthVal_S[i]
9:     if diff < min then
10:       min = diff
11:       idx ← i
12:     end if
13:   end for
14:   preFinalDepthVal =  $\frac{\text{depthVal\_S}[\text{idx} + \frac{n}{2}] + \text{depthVal\_S}[\text{idx}]}{2}$ 
15: end procedure

```

327 computation of the minimum of medians m (stated as
 328 idx in Algorithm 1). Here, before beginning the calcu-
 329 lation of obtaining m and corresponding depth value, we
 330 sort the depth values d_i (see line 5 of Algorithm 1) and
 331 split the traversal path to half of the window size. The
 332 validity and efficiency of using half of the sorted sam-
 333 ple data (i.e. splitting the sorted depth values into half)
 334 rather than checking each data to the rest of the other
 335 data one by one (which we previously did in Figure 7)
 336 is stated in [20]. We adapt this time-efficient median
 337 calculation process into 1DLMedS to reduce the overall
 338 residual (i.e. difference between the depth values) cal-
 339 culation time for each of the sorted depth values inside
 340 the window. After we obtain the minimum difference
 341 (see line 10 of Algorithm 1), we find the index idx of
 342 the minimum difference (see line 11) and retrieve the
 343 depth value which has this minimum index idx (see line
 344 14). Based on this depth value and index idx , we calcu-
 345 late the $\hat{\sigma}$, and find the inliers and outliers according to
 346 w_i of Equation 4. Then we take the average of the depth
 347 values of the inliers as the final depth value.

348 In this way, 1D LMedS enhances a certain depth
 349 frame t_i by assigning stable valid depth values for each
 350 pixel of t_i . Then the sliding window moves one frame
 351 further to enhance the next frame t_{i+1} . It is worth to
 352 mention that, from the $(n+1)$ th frame ($(t+4)$ th frame
 353 on the output stream in Figure 4) we would observe the
 354 enhanced depth frame. For the $(n+2)$ th frame ($(t+5)$ th
 355 frame in Figure 4) and later ones, we just have to wait
 356 for one extra frame's calculation.

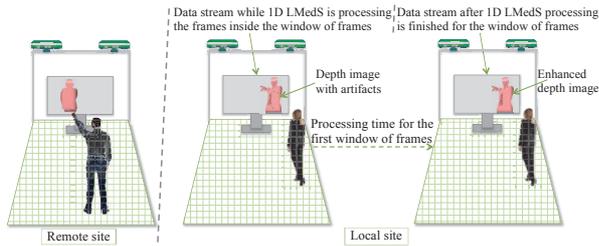


Figure 8: Illustration of using 1D LMedS in a telepresence scenario – the depth image on the right screen of the local site is enhanced by 1D LMedS while the left depth image shows artifacts (since the processing of the frames inside the sliding window is in progress).

357 6. Example use case: telepresence

358 1D LMedS can be applied to enhance the depth
 359 frames in many application areas such as in 3D re-
 360 construction of a scene, tracking or forensic analysis.
 361 Telepresence, an emerging application area which often
 362 needs 3D scene reconstruction, can benefit from the 1D
 363 LMedS depth enhancement method. In case of telepres-
 364 ence, our 1D LMedS can be used to enhance the stream
 365 of frames received from the remote site of the commu-
 366 nication. There are typically two sides of communica-
 367 tion in a traditional telepresence system, which are: lo-
 368 cal site and remote site. Users staying at each site can
 369 communicate and interact with each other via network
 370 channel; their surrounding environment is captured and
 371 transmitted to the other end of the communication site.
 372 Figure 8 illustrates a typical telepresence scenario along
 373 with how 1D LMedS can be applied to such scenario.

374 Since, in case of 1D LMedS, we use a sliding win-
 375 dows of frames to enhance the t_i th frame, we will see the
 376 enhanced t_i th frame after processing the frames inside
 377 the sliding window. While the frames inside the sliding
 378 window are being processed, the users will observe the
 379 frames with artifacts, but as soon as the processing is
 380 finished and the sliding window moves to the next win-
 381 dows of frames, they will start observing the enhanced
 382 stream. As stated in the last paragraph of Section 5, for
 383 the later sliding window, only one frame’s worth calcu-
 384 lation time is needed since the 1D LMedS calculation
 385 for others frames have already been done while process-
 386 ing the frames inside the previous sliding window. Simi-
 387 larly, 1D LMedS can be applied to other applications
 388 (e.g. in tracking) which involves a sequence of frames
 389 to be processed before carrying out a certain task.

390 7. Results and discussion

391 In this section, we perform a set of experiments
 392 to evaluate the overall performance of proposed 1D

393 LMedS approach of enhancing a depth image sequence.
 394 Here, we use several datasets among which some are
 395 captured locally with Kinect v1 and three others (one is
 396 from Camplani et. al. [8] and two are from Middlebury
 397 RGB-D database [39]) are used as reference datasets.
 398 Here we use depth data from Kinect v1 sensor because,
 399 most of the reference works have been pursued by using
 400 Kinect v1 sensor. The local datasets and the dataset
 401 from [8] have the resolution of 480×640 while Mid-
 402 dlebury datasets have 555×690 resolution. We use the
 403 Middlebury datasets for the quantitative evaluation and
 404 other datasets for the qualitative evaluation.

405 To compare the results, we use five state-of-the-art
 406 methods which are the joint bilateral filtering approach
 407 [8, 28] from Camplani et. al., the depth-color fusion
 408 method [9] from Garcia et. al., depth smoothing method
 409 [37] from Nguyen et. al. and the trilateral approach
 410 [31] from Wang et.al. To achieve a fair comparison sce-
 411 nario during the experiments, we tuned the parameters
 412 of these methods to achieve the best results as stated in
 413 those respective methods. Our CPU implementation of
 414 1D LMedS is performed on a Macbook Pro with an Intel
 415 Core i7 2.3 GHz processor and 8GB RAM. In Section
 416 7.1 and Section 7.2, we present the qualitative and the
 417 quantitative performance comparison respectively.

418 7.1. Performance on real-world depth sequences

419 We tested 1D LMedS on various datasets which are
 420 captured locally and on a dataset from [8]. Since we did
 421 not find any reference sequence of depth frames which
 422 was captured with a camera that was kept static while
 423 recording, we used local datasets and image sequences
 424 from [8] to evaluate the performance of 1D LMedS.

425 Figure 9 shows the performance of 1D LMedS on
 426 three different datasets. Here, the images on the top
 427 three rows are captured locally and the images of 4th
 428 row are taken from [8]. Here, the images in the 2nd
 429 row and the images in the 3rd, 4th and 5th column of
 430 3rd and 4th rows of Figure 9 show the performance of
 431 1D LMedS with different number of frames in the slid-
 432 ing window. As expected, with the number of frames
 433 increased in the sliding window, the chance of getting
 434 better values for the pixels are also increased and hence,
 435 we observe that more and more artifacts are removed
 436 with bigger sliding windows. In Figure 10 we show the
 437 comparison results among two reference depth enhance-
 438 ment methods ([9] and [37]) and our 1D LMedS on local
 439 datasets and the dataset of [8]. From this figure, we can
 440 see that 1D LMedS performs quite well in removing the
 441 artifacts compared to the reference methods.

442 We applied 1D LMedS to both static and moving
 443 parts of the captured scene. 1D LMedS worked quite

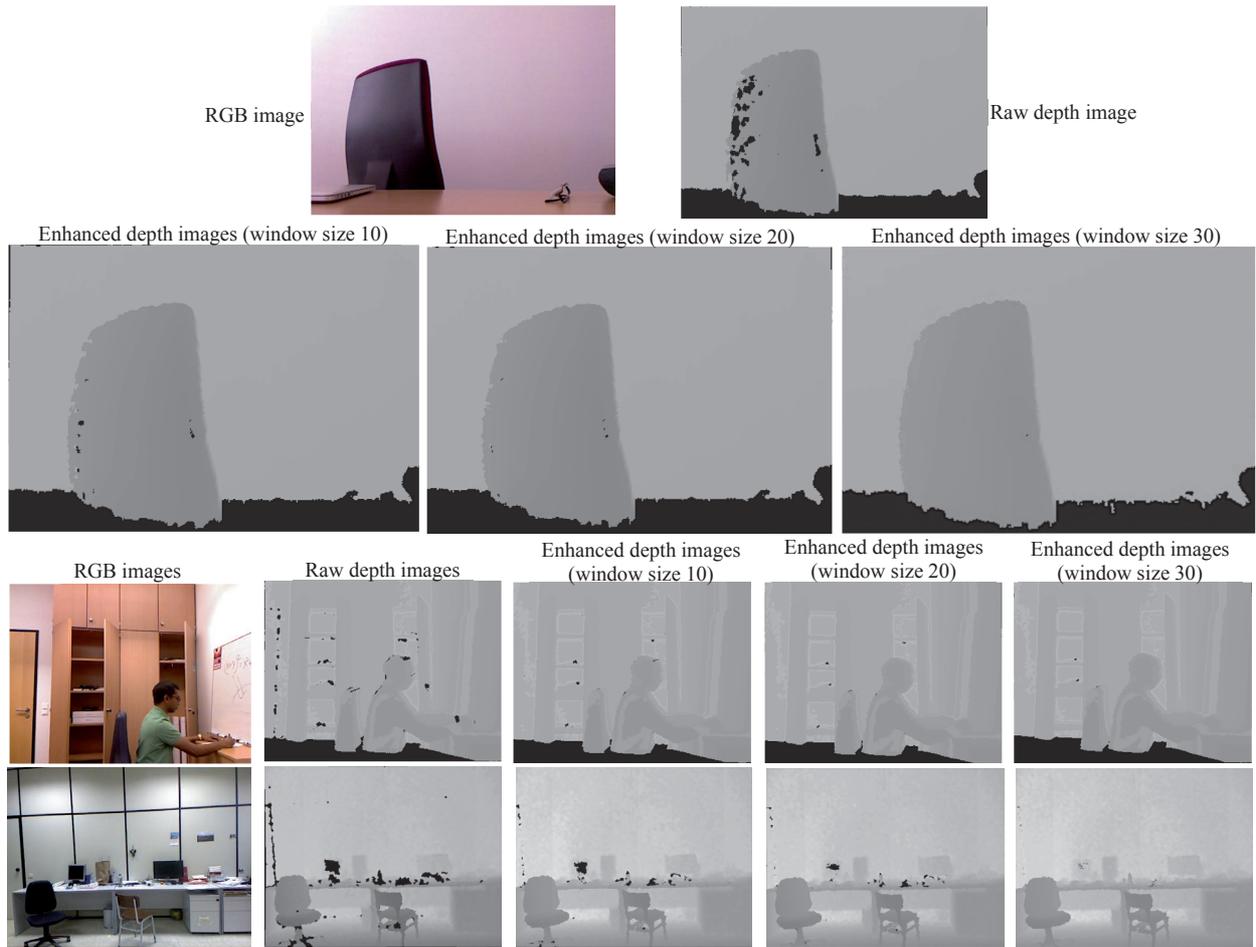


Figure 9: 1D LMedS performance with varying window size on different datasets – images on the 2nd row and on the 3rd, 4th and 5th columns of 3rd and 4th rows are the enhanced depth images by 1D LMedS with window size 10, 20 and 30 respectively. The test image on the 4th row is the 20th frame which is taken from dataset [8] and the images on the other three rows are captured locally on various days.

444 well with static parts of the scene (see Figure 11(a –
 445 e)) while stabilizing the depth values over consecutive
 446 frames (see Figure 11(d – e)). We can also perceive the
 447 performance of 1D LMedS by comparing Figure 11(e)
 448 to Figure 2(b). Here, we used 25 frames in the sliding
 449 window. For the moving parts of the scene, 1D LMedS
 450 also performed good in filling the missing depth val-
 451 ues (see Figure 12 and Figure 13). It is worth to men-
 452 tion that, for a scene with moving objects, we use a
 453 smaller sliding window size than a scene with static ob-
 454 jects while the implementation of 1D LMedS remains
 455 the same in both cases. If we use similar window sizes
 456 in dynamic scenes as in static scenes, we would ob-
 457 serve ghosting artifacts in the resulting images (see the
 458 gray shadow near the edge of the chair, marked with
 459 purple rectangle and zoomed for clear visualization, in
 460 Figure 13). Another issue to mention here is, when a

461 pixel does not have a valid depth value throughout the
 462 whole sliding window, that pixel continues to remain as
 463 a black spot. It is also true for a large area with missing
 464 depth value (see the orange ellipse marked area in Fig-
 465 ure 13). Increasing window size in the hope for finding
 466 valid depth values for such pixels might work but that
 467 would result in lower computation speed and increased
 468 overall latency [40]. In Figure 12, we show the original
 469 and enhanced images after pixel aligning the original
 470 and enhanced depth images to their corresponding color
 471 images; the images are rendered in OpenGL. For the re-
 472 sults of Figure 12 and 13, we used 30 and 35 frames,
 473 respectively, in the sliding window.

474 With a bit more careful observation of the enhanced
 475 depth images, we found out that our 1D LMedS pre-
 476 serves the sharpness of the edges of different objects
 477 better than most of these mentioned methods, see the red

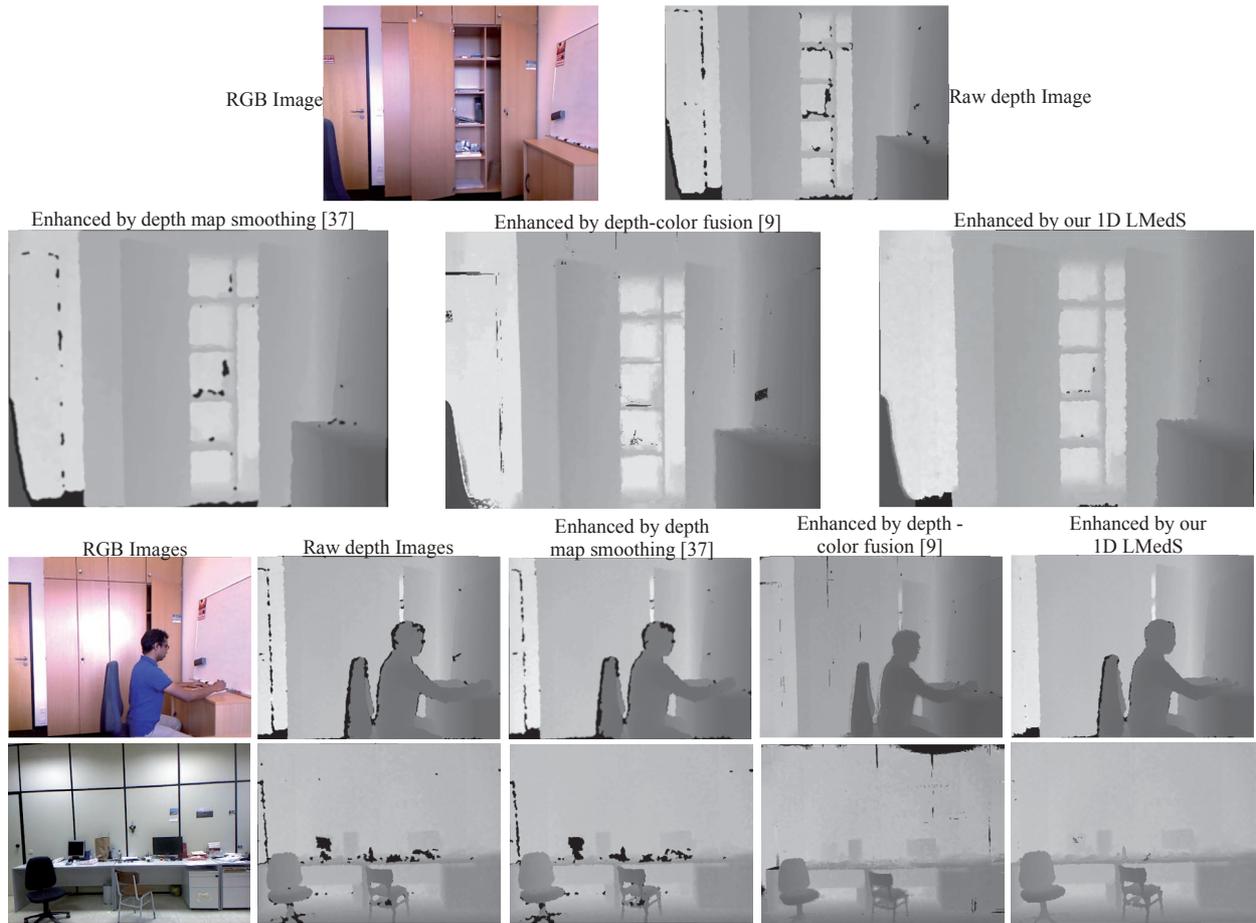


Figure 10: Comparison of results on datasets captured locally and from [8] with our 1D LMedS, depth smoothing method [37] and fusion based enhancement[9] – the images on 2nd row and on the 3rd, 4th and 5th columns of 3rd and 4th rows are the enhanced depth images by method [37], by method [9] and by our 1D LMedS respectively. The test image on the 4th row is the 20th frame of the dataset [8] and the images on the other three rows are captured locally.

478 marked area and the zoomed parts in Figure 14. Here,
 479 we take the method of [31] to compare the enhancement
 480 results. On the surface of the objects, the method from
 481 [31] produces quite nice results, however it suffers along
 482 the edge of the objects. Moreover, we also notice that
 483 the method in [31] highly depends on a very accurate
 484 registration between the color and the depth image pairs.
 485 If the color image and depth image pair is not perfectly
 486 aligned, this method shows blurry object edges. This is
 487 due to the fact that the color of an object in the color
 488 image can be considered to be belonging to a different
 489 object in the corresponding depth image pair. Here, for
 490 example on the first and third row in Figure 14, the edge
 491 of the nose and mouth of the person, and the edge of the
 492 right side of the chair are mistakenly considered to be
 493 belonging to both the background and the foreground.
 494 This results in blurry edges which have depth values

495 that are mixed from depth values of different objects.
 496 Besides this, when the background and the foreground
 497 objects have same color (see the diamond marked areas
 498 on the color and depth images on the first row of Fig-
 499 ure 14; here the hair of the person and the background
 500 have the same color), it also produces artifacts near the
 501 edges of the objects (see the yellow dotted region on the
 502 zoomed part of Figure 14) in the depth image. Here too,
 503 the algorithm mixes the depth values from both objects.
 504 In 1D LMedS, we do not face these issues since we only
 505 need the depth sequence for the enhancement.

506 Moreover, 1D LMedS takes far less time to enhance
 507 a single frame than other methods. For example, to en-
 508 hance one frame of Dataset 1, 1D LMedS takes only
 509 0.14 second (see the value in the 4th column of 2nd row
 510 in Table 1), whereas, the method in [31] takes more
 511 than 20 minutes (1200 seconds) due to its time consum-

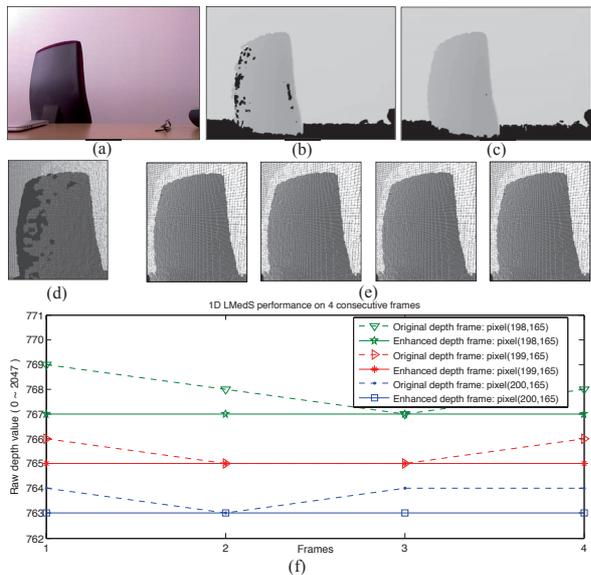


Figure 11: **Results with static object** – (a) color image, (b,d) raw depth image, (c) enhanced depth image by 1D LMedS, (e) 4 consecutive depth frames (zoomed for clear visibility) with stable depth values, (f) the depth values (before and after applying 1D LMedS) of 3 pixels on the 4 frames of (b) show that the enhanced frames have stable depth values for all of these 3 pixels over the consecutive frames.



Figure 12: **Results-1 with moving object** – (a) image with raw depth map, (b) image with enhanced depth map by 1D LMedS. Here, the depth images are pixel aligned to their corresponding color images.

512 ing energy minimization computation. Due to this high
 513 computation, similar color image based methods such
 514 as [9, 30–32] are not suitable for real-time applications.
 515 Table 1 shows the performance of proposed 1D
 516 LMedS in terms of frames per second (*fps*) on the four
 517 datasets we used in our experiments. In Table 1, we
 518 refer *Dataset 1* to the images on the 1st and 2nd rows of
 519 Figure 9, *Dataset 2* to the images on the 3rd row of Fig-
 520 ure 9, *Dataset 3* to the images on the 4th row of Figure 9
 521 and *Dataset 4* to the images on the 2nd row of Figure 10.
 522 We show here the performance by varying the number
 523 of frames in the sliding window. As an example, for the
 524 *Dataset 1* with 25 frames in the sliding window ($ws =$
 525 25), we achieved a performance of 8 *fps* with our CPU
 526 implementation of 1D LMedS.

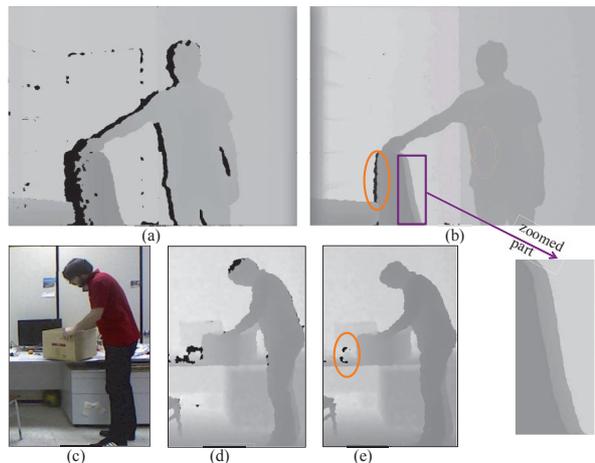


Figure 13: **Results-2 with moving object** – (a,d) raw depth images, (b,e) enhanced depth images by 1D LMedS, (c) color image (part of frame 244 from [8]). The black area marked with orange ellipses in (b) and (e) occurs because the depth sensor cannot visualize that area and hence, there is no depth information for that area. The ghosting artifact (a trail of gray shadow) marked with purple rectangle and zoomed in (b) occurs due to two reasons: the object is moving and we take 35 frames in the sliding window.

Table 1: Performance, in frames per second (*fps*), of 1D LMedS by varying the sliding window size (*ws*) on four local datasets.

Datasets	<i>fps</i> (<i>ws</i> :10)	<i>fps</i> (<i>ws</i> :20)	<i>fps</i> (<i>ws</i> :25)	<i>fps</i> (<i>ws</i> :30)
Dataset 1	12	10	8	7
Dataset 2	13	11	8	7
Dataset 3	12	9	7	6
Dataset 4	12	10	7	6

527 In Figure 15 and Figure 16, we depict a comparison
 528 result between the method [37] and 1D LMedS on a se-
 529 quence of 100 frames on two local datasets. The depth
 530 values plotted on both figures are for the pixels marked
 531 by the yellow dot on the depth images. In both cases, we
 532 observe that our proposed 1D LMedS outperforms the
 533 method [37] in removing the invalid depth values while
 534 stabilizing the valid depth values.

535 7.2. Performance on artificially degraded Kinect depth 536 sequences

537 Here, we perform yet another comparison experiment
 538 for an extensive evaluation of the performance of our
 539 1D LMedS. This time we take reference datasets (the
 540 'Art' and the 'Book') which is publicly available in the
 541 Middlebury depth datasets [39]. Since we perform a
 542 quantitative experiment here, we use standard metrics
 543 such as PSNR (Peak Signal-to-Noise Ratio) and SSIM

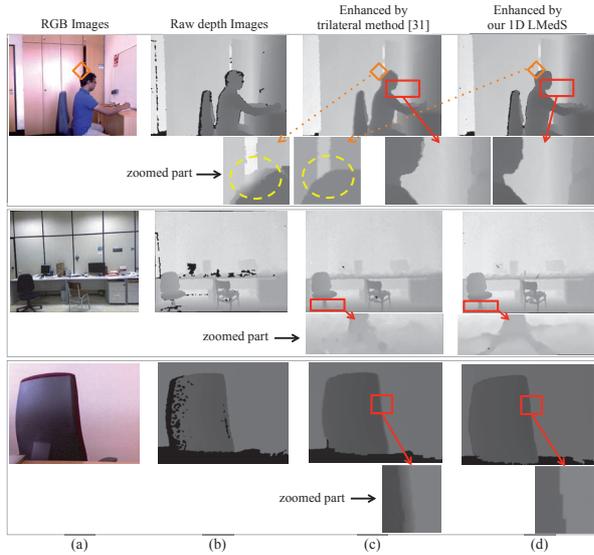


Figure 14: Detailed comparison of results based on the marked areas of the three datasets with our 1D LMedS and trilateral method [31] – (a) color images, (b) depth images, (c, d): enhanced depth images by trilateral method [31], by our 1D LMedS method respectively. The test image on the 2nd row is the 20th frame of the dataset [8] and the images on the other 2 rows are captured locally.

544 (Structural SIMilarity) index for the reference datasets.
 545 Now, since ground truth depth images are required to
 546 compute the PSNR and SSIM metrics and, since Mid-
 547 lebury datasets provide the ground truth depth images,
 548 we follow the procedure stated in [32] to generate arti-
 549 ficially Kinect-like degraded depth images. In [32],
 550 the authors create the structural missing phenomenon
 551 of Kinect’s depth images along the areas with depth dis-
 552 continuities and the random missing phenomenon on
 553 the areas which are flat.

554 Since we need a sequence of frames for 1D LMedS
 555 and since the Middlebury datasets do not provide se-
 556 quence of depth images, we take multiple copies of
 557 ‘Art’ and ‘Book’ depth images and then apply the
 558 artificial degradation process from [32] randomly on
 559 each of these images. Hence, we get a sequence of
 560 frames which contain Kinect-like degradation on dif-
 561 ferent places for each of the frames. Then, we apply
 562 1D LMedS on these sequence of frames and get the
 563 final enhanced depth image. We can consider this se-
 564 quence of frames as a part of a static scene where arti-
 565 facts occurs over the temporal domain, which we have
 566 also seen in our previous experiment in Figure 11. For
 567 the other methods [28, 31] to which we compare the
 568 performance of 1D LMedS, we use the color and depth
 569 images from the ‘Art’ and ‘Book’ datasets, and the arti-
 570 ficially Kinect-like degraded image from [32]. For these

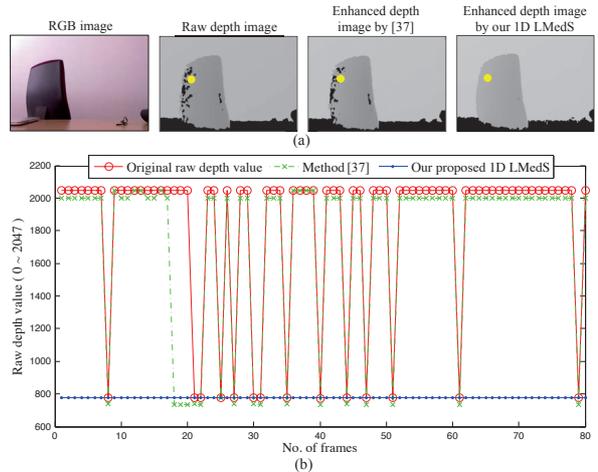


Figure 15: Comparison of results-1 with method [37] and proposed 1D LMedS on 100 frames of a local dataset – (a) from left to right – color image, raw depth image, enhanced depth images by [20] and by 1D LMedS, (b) depth enhancement results (for the pixel marked by yellow circle) by method [37] and 1D LMedS. 1D LMedS removes more invalid depth values than method [37] and it also stabilizes the valid depth values which is shown by the blue filled-circle marked line in (b). In (a), the enhanced frame is the 27th frame of the 100 frames and in (b), we show the depth values for 80 enhanced frames (sliding window contains the first 20 frames).

Table 2: Quantitative comparison results of different methods on standard datasets with artificial Kinect-like degradations. Here, for our 1D LMedS, sliding window size (ws) is set to 20.

	Datasets			
	Art		Book	
	PSNR (dB)	SSIM	PSNR (dB)	SSIM
Camplani [28]	29.1845	0.9442	26.5911	0.9416
Wang [31]	30.3104	0.9587	26.8854	0.9492
Our 1D LMedS	32.2324	0.9723	28.3123	0.9624

571 methods [28, 31], the depth recovery process is per-
 572 formed on this artificially degraded depth image.

573 To compute the PSNR and SSIM metrics, we con-
 574 sider the ground truth depth images from the ‘Art’ and
 575 ‘Book’ datasets as the reference depth images, and the
 576 enhanced depth images as the images to be compared
 577 against the reference depth images. Figure 17 shows
 578 the color images, artificially degraded depth images and
 579 enhanced depth images results by [28], [31] and our
 580 1D LMedS respectively. Table 2 shows the PSNR and
 581 SSIM metrics which reflects the quantitative compari-
 582 son among these methods.

583 When we take a closer look at the results (see Fig-
 584 ure 17) of different depth image enhancement methods,
 585 we observe that all of them perform fine to recover the
 586 missing depth values on the surface areas which are flat.
 587 However, the difference among these methods lies in

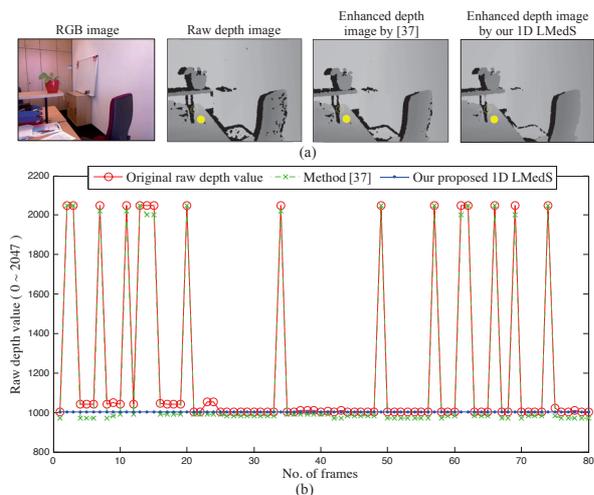


Figure 16: Comparison of results-2 with method [37] and 1D LMedS on 100 frames of another local dataset – here also, 1D LMedS removes more invalid depth values than method [37] and stabilizes the valid depth values (see the blue filled-circled line in (b)).

Table 3: Quantitative comparison results of 1D LMedS by varying the sliding window size (ws) on standard datasets with artificial Kinect-like degradations.

Window size(ws)	Datasets			
	Art		Book	
	PSNR (dB)	SSIM	PSNR (dB)	SSIM
10	31.7458	0.9651	27.8957	0.9573
20	32.2324	0.9723	28.3123	0.9624
25	32.6981	0.9774	28.7923	0.9682
30	33.0148	0.9802	29.1635	0.9710

the enhanced areas near the edges of the objects where the structural discontinuity occurred inside the degraded depth images. Although, all of these methods recover the edges, the results from methods [28] and [31] show blurry object edges, whereas our 1D LMedS preserves the sharpness of the object’s edges. The result of this subjective comparison is also reflected in the PSNR and SSIM metrics in Table 2 where 1D LMedS outperforms the other methods. We show the PSNR and SSIM metrics of our 1D LMedS with different sliding window size (ws) on the Middlebury datasets with artificial Kinect-like degradations in Table 3.

As an alternative of using the depth values of only the current pixel in the sliding window of frames, we also considered using regional information surrounding the current pixel. We can use a patch of size $p \times p$ and take the average of the depth values of the pixels inside this patch and use that value as the current pixel’s depth value while processing the frames in the sliding window. Here, the method of 1D LMedS remains the

same, except, instead of taking just one depth value from the current pixel, we consider the patch as a vector of consecutive depth values. We test this approach for the pixels with invalid depth values. While testing this idea, we found out that taking a smaller patch (3×3) and smaller sliding window size doesn’t necessarily improve the depth enhancement performance; because, there were not enough valid depth values in the patches. However, we got a slightly better result with patch size 7×7 ; the improvements were mostly in the areas which have large areas with invalid depth values. However, we observe a noticeable decrease in computational speed and increase in latency with such bigger patch size. Therefore, we decide not to include the outcome from this idea in this article and plan to improve on this idea with a better formulation and keep it as future work.

7.3. Limitations

Our approach has some limitations. Firstly, we cannot recover invalid depth values or noise that persist in one region of the depth image sequence with our 1D LMedS method (see the black area marked with orange ellipses in Figure 13(b) and (e)). In this case, regional information, such as using a patch of neighboring pixels, might be useful to refine the large areas with artifacts. Secondly, we may see ghosting artifacts if we set the size of the sliding window too large (see the gray-shadowed area marked with a purple rectangle in Figure 13(b)). Ghosting artifacts may also occur for fast-moving objects, where depth values change drastically from one frame to another. We considered using our previous method from [35] based on optical flow motion vectors to further reduce ghosting. That method, however, works well with small temporal windows sizes (the experiments in [35] used a window with five image frames) and increasing the window size might increase ghosting due to accumulating errors in optical flow calculations. Thus, optimal temporal windows sizes for the method described in this paper and the method from [35] are not compatible. For future work, we can envision a meta algorithm that classifies image regions based on optical flow. The method from [35] could then be applied to regions with large optical flow, while the method proposed in this paper could be applied to regions with smaller optical flow. Although, 1D LMedS exhibits fewer edge-bleeding artifacts, it could be further reduced by using an anisotropic diffusion approach along with structural similarity of depth and color frames. Moreover, a GPU implementation of 1D LMedS would allow us to enhance the depth frames in real time.

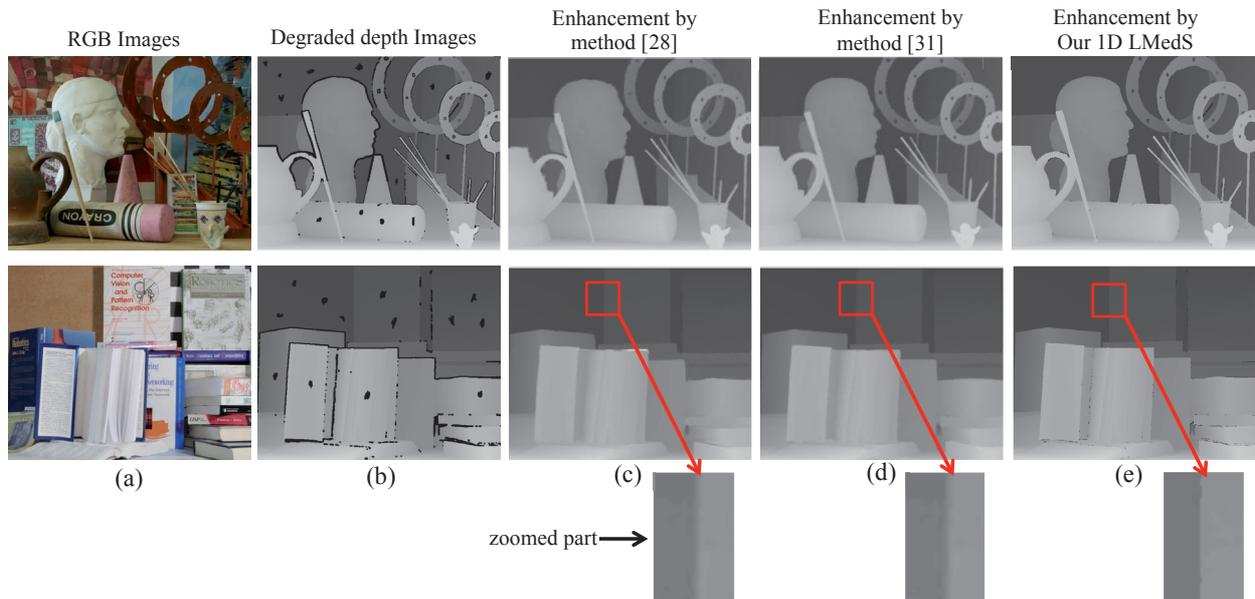


Figure 17: Subjective comparison results on datasets 'Art' and 'Book' from Middlebury database [39] with artificial Kinect-like degradations – a) color images, (b) degraded depth images, (c – e): enhanced depth images by method [28], method [31] and our 1D LMedS respectively. The zoomed parts below the 2nd row show detailed comparison, among these methods, on preserving the sharpness of the object's edges.

659 8. Conclusion

660 We have introduced a robust 1D LMedS approach
 661 which can efficiently enhance the depth images gener-
 662 ated by depth cameras. The output of our proposed ap-
 663 proach shows good performance in filling the invalid
 664 depth values and stabilizing the valid values for static
 665 and dynamic parts of a scene. We have used local
 666 datasets and reference datasets to show the performance
 667 of 1D LMedS. For all the test datasets, our 1D LMedS
 668 method exhibits good performance. Moreover, we have
 669 also used other reference methods on the test datasets
 670 to compare results with 1D LMedS; in this case also,
 671 1D LMedS performs quite good comparing with other
 672 methods. The proposed 1D LMedS is computationally
 673 less expensive than higher dimensional LMedS and its
 674 implementation procedure on depth frames is efficient.

675 Since the proposed 1D LMedS performs quite well
 676 in removing the invalid depth values and stabilizing the
 677 inconsistent valid depth values (considered as outliers),
 678 one might apply 1D LMedS to address the so called 'fly-
 679 ing pixel' [41] issues found in the content generated by
 680 ToF sensors. Such pixel's depth values could also be
 681 identified as outliers and be stabilized with 1D LMedS.

682 Acknowledgements

683 This work was supported by the EU FP7 Marie Curie
 684 ITN "DIVA" under REA Grant Agreement No.290227

685 and the DFG Research Grant STA 708/5-1 "SMOOTH".

686 We would like to thank Dr. Massimo Camplani (Grupo
 687 de Tratamiento de Imágenes, Universidad Politécnica
 688 de Madrid) for providing dataset and Dr. Frederic Gar-
 689 cia (Interdisciplinary Centre for Security, Reliability
 690 and Trust, University of Luxembourg) for his valuable
 691 comments during the preparation of this article.

692 References

- 693 [1] A. Maimone, H. Fuchs, Encumbrance-free telepresence system
 694 with real-time 3d capture and display using commodity depth
 695 cameras, in: 10th IEEE ISMAR, 2011, pp. 137–146.
 696 [2] M. Urschler, J. Hiller, A. Bornik, T. Paul, M. Giretzlehner,
 697 H. Bischof, K. Yen, E. Scheurer, Intuitive presentation of clinical
 698 forensic data using anonymous and person-specific 3d refer-
 699 ence manikins, *Forensic Science International* 241 (0) (2014)
 700 155 – 166.
 701 [3] M. Dou, H. Fuchs, Temporally enhanced 3d capture of room-
 702 sized dynamic scenes with commodity depth cameras, in: *IEEE*
 703 *Virtual Reality (VR)*, 2014, pp. 39–44.
 704 [4] S. Matyunin, D. Vatolin, Y. Berdnikov, M. Smirnov, Temporal
 705 filtering for depth maps generated by kinect depth camera, in:
 706 *3DTV Conference: The True Vision - Capture, Transmission*
 707 *and Display of 3D Video (3DTV-CON)*, 2011, pp. 1–4.
 708 [5] A. Amamra, N. Aouf, Gpu-based real-time rgbd data filtering,
 709 *Journal of Real-Time Image Processing* (2014) 1–18.
 710 [6] K. He, J. Sun, X. Tang, Guided image filtering, in: *11th Euro-*
 711 *pean Conference on Computer Vision (ECCV)*, Springer Berlin
 712 Heidelberg, 2010, pp. 1–14.
 713 [7] J. Liu, X. Gong, J. Liu, Guided inpainting and filtering for kinect
 714 depth maps, in: *21st International Conference on Pattern Recog-*
 715 *nition (ICPR)*, 2012, pp. 2055–2058.

- 716 [8] M. Camplani, T. Mantecon, L. Salgado, Depth-color fusion
717 strategy for 3-d scene modeling with kinect, *IEEE Transactions*
718 *on Cybernetics* 43 (6) (2013) 1560–1571.
- 719 [9] F. Garcia, D. Aouada, T. Solignac, B. Mirbach, B. Ottersten,
720 Real-time depth enhancement by fusion for rgb-d cameras,
721 *Computer Vision, IET* 7 (5) (2013) 1–11.
- 722 [10] S. Beck, A. Kunert, A. Kulik, B. Froehlich, Immersive group-
723 to-group telepresence, *IEEE Transactions on Visualization and*
724 *Computer Graphics* 19 (4) (2013) 616–625.
- 725 [11] L. Chen, H. Lin, S. Li, Depth image enhancement for kinect
726 using region growing and bilateral filter, in: *21st International*
727 *Conference on Pattern Recognition*, 2012, pp. 3070–3073.
- 728 [12] Q. Yang, N. Ahuja, R. Yang, K.-H. Tan, J. Davis, B. Culbertson,
729 J. Apostolopoulos, G. Wang, Fusion of median and bilateral filter-
730 ing for range image upsampling, *IEEE Transactions on Image*
731 *Processing* 22 (12) (2013) 4841–4852.
- 732 [13] J. Dolson, J. Baek, C. Plagemann, S. Thrun, Upsampling range
733 data in dynamic environments, in: *IEEE Computer Vision and*
734 *Pattern Recognition (CVPR)*, 2010.
- 735 [14] F. Qi, J. Han, P. Wang, G. Shi, F. Li, Structure guided fusion for
736 depth map inpainting, *Pattern Recognition Letter* 34 (1) (2013)
737 70–76.
- 738 [15] C. Richardt, C. Stoll, N. A. Dodgson, H.-P. Seidel, C. Theobalt,
739 Coherent spatiotemporal filtering, upsampling and rendering of
740 rgbz videos, *Computer Graphics Forum* 31 (2) (2012) 247–256.
- 741 [16] B.-S. Lin, M.-J. Su, P.-H. Cheng, P.-J. Tseng, S.-J. Chen, Tem-
742 poral and spatial denoising of depth maps, *Sensors* 15 (8) (2015)
743 18506–18525.
- 744 [17] S.-Y. Kim, J.-H. Cho, A. Koschan, M. Abidi, Spatial and tem-
745 poral enhancement of depth images captured by a time-of-
746 flight depth sensor, in: *20th International Conference on Pattern*
747 *Recognition (ICPR)*, 2010, pp. 2358–2361.
- 748 [18] A. B. M. T. Islam, C. Scheel, R. Pajarola, O. Staadt, Robust
749 enhancement of depth images from kinect sensor, in: *2015 IEEE*
750 *Virtual Reality (VR)*, 2015, pp. 197–198.
- 751 [19] A. B. M. T. Islam, C. Scheel, R. Pajarola, O. Staadt, Depth im-
752 age enhancement using 1d least median of squares, in: *Com-*
753 *puter Graphics International (CGI'15)*, 2015.
- 754 [20] P. J. Rousseeuw, Least median of squares regression, *Journal of*
755 *the American statistical association* 79 (388) (1984) 871–880.
- 756 [21] S. Bhattacharya, S. Gupta, K. Venkatesh, High accuracy depth
757 filtering for kinect using edge guided inpainting, in: *International*
758 *Conference on Advances in Computing, Communications and*
759 *Informatics (ICACCI)*, 2014, pp. 868–874.
- 760 [22] L. Zhang, P. Shen, S. Zhang, J. Song, G. Zhu, Depth enhance-
761 ment with improved exemplar-based inpainting and joint trilat-
762 eral guided filtering, in: *IEEE International Conference on Im-*
763 *age Processing (ICIP)*, 2016, pp. 4102–4106.
- 764 [23] T.-W. Hui, K. N. Ngan, Motion-depth: Rgb-d depth map en-
765 hancement with motion and depth in complement, in: *IEEE*
766 *Conference on Computer Vision and Pattern Recognition*
767 *(CVPR)*, 2014, pp. 3962–3969.
- 768 [24] K. Lai, L. Bo, X. Ren, D. Fox, A large-scale hierarchical multi-
769 view rgb-d object dataset, in: *IEEE International Conference on*
770 *Robotics and Automation (ICRA)*, 2011, pp. 1817–1824.
- 771 [25] J. Wasza, S. Bauer, J. Hornegger, Real-time preprocessing for
772 dense 3-d range imaging on the gpu: Defect interpolation, bi-
773 lateral temporal averaging and guided filtering, in: *International*
774 *Conference on Computer Vision Workshops*, 2011, pp. 1221–
775 1227.
- 776 [26] A. Maimone, J. Bidwell, K. Peng, H. Fuchs, Enhanced personal
777 autostereoscopic telepresence system using commodity depth
778 cameras, *Computers & Graphics* 36 (7) (2012) 791 – 807.
- 779 [27] A. Telea, An image inpainting technique based on the fast
780 marching method, *Graphics Tools* 9 (1) (2004) 25–36.
- 781 [28] M. Camplani, L. Salgado, Efficient spatio-temporal hole filling
782 strategy for kinect depth maps, in: *Proc. SPIE*, Vol. 8290, 2012,
783 pp. 82900E–82900E–10.
- 784 [29] J. Zhu, L. Wang, R. Yang, J. Davis, Fusion of time-of-flight
785 depth and stereo for high accuracy depth maps, in: *IEEE Con-*
786 *ference on Computer Vision and Pattern Recognition*, 2008, pp.
787 1–8.
- 788 [30] L. Sheng, K. N. Ngan, Depth enhancement based on hybrid ge-
789 ometric hole filling strategy, in: *2013 IEEE International Con-*
790 *ference on Image Processing*, 2013, pp. 2173–2176.
- 791 [31] Z. Wang, J. Hu, S. Wang, T. Lu, Trilateral constrained sparse
792 representation for kinect depth hole filling, *Pattern Recognition*
793 *Letter* 65 (C) (2015) 95–102.
- 794 [32] J. Yang, X. Ye, K. Li, C. Hou, Y. Wang, Color-guided depth re-
795 covery from rgb-d data using an adaptive autoregressive model,
796 *IEEE Transactions on Image Processing* 23 (8) (2014) 3443–
797 3458.
- 798 [33] S. Milani, G. Calvagno, Correction and interpolation of depth
799 maps from structured light infrared sensors, *Signal Processing: Image*
800 *Communication* 41 (2016) 28 – 39.
- 801 [34] Y. Yang, Q. Liu, R. Ji, Y. Gao, Dynamic 3d scene depth recon-
802 struction via optical flow field rectification, *Public Library of*
803 *Science (PloS ONE)* 7 (11) (2012) e47041.
- 804 [35] R. Avetisyan, C. Rosenke, M. Luboschik, O. Staadt, Tem-
805 poral filtering of depth images using optical flow, in: *Pro-*
806 *ceedings of the 24th International Conference in Central Eu-*
807 *rope on Computer Graphics, Visualization and Computer Vision*
808 *(WSCG'16)*, 2016.
- 809 [36] K. Khoshelham, E. O. Elberink, Accuracy and resolution of
810 kinect depth data for indoor mapping applications, *Sensors*
811 12 (2) (2012) 1437–1454.
- 812 [37] C. Nguyen, S. Izadi, D. Lovell, Modeling kinect sensor noise for
813 improved 3d reconstruction and tracking, in: *2nd International*
814 *Conference on 3D Imaging, Modeling, Processing, Visualiza-*
815 *tion and Transmission (3DIMPVT)*, 2012, pp. 524–530.
- 816 [38] P. J. Rousseeuw, A. M. Leroy, *Robust Regression and Outlier*
817 *Detection*, John Wiley & Sons, Inc., New York, NY, USA, 1987.
- 818 [39] MiddleburyDatasets, [http://vision.middlebury.edu/](http://vision.middlebury.edu/stereo/data/)
819 [stereo/data/](http://vision.middlebury.edu/stereo/data/), accessed: 2017-05-16 (2013).
- 820 [40] S. Ohl, M. Willert, O. Staadt, Latency in distributed acquisition
821 and rendering for telepresence systems, *IEEE Tran. on Visual-*
822 *ization and Computer Graphics* 21 (12) (2015) 1442–1448.
- 823 [41] A. Kolb, E. Barth, R. Koch, R. Larsen, Time-of-flight sensors
824 in computer graphics, in: *Proc. Eurographics (State-of-the-Art*
825 *Report)*, 2009.