

Local self-similarity based registration of human ROIs in pairs of stereo thermal-visible videos

Atousa Torabi^{*,a}, Guillaume-Alexandre Bilodeau^a

^a*LITIV, Department of Computer and Software Engineering,
École Polytechnique de Montréal,
P.O. Box 6079, Station Centre-ville, Montréal
(Québec), Canada, H3C 3A7*

Abstract

For several years, mutual information (MI) has been the classic multimodal similarity measure. The robustness of MI is closely restricted by the choice of MI window sizes. For unsupervised human monitoring applications, obtaining appropriate MI window sizes for computing MI in videos with multiple people in different sizes and different levels of occlusion is problematic. In this work, we apply local self-similarity (LSS) as a dense multimodal similarity metric and show its adequacy and strengths compared to MI for a human ROIs registration. We also propose a LSS-based registration of thermal-visible stereo videos that addresses the problem of multiple people and occlusions in the scene. Our method improves the accuracy of the state-of-the-art disparity voting (DV) correspondence algorithm by proposing a motion segmentation step that approximates depth segments in an image and enables assigning disparity to each depth segment using larger matching

*Corresponding author

Email addresses: atousa.torabi@polymtl.ca (Atousa Torabi),
guillaume-alexandre.bilodeau@polymtl.ca (Guillaume-Alexandre Bilodeau)

window while keeping registration accuracy. We demonstrate that our registration method outperforms the recent state-of-the-art MI-based stereo registration for several realistic close-range indoor thermal-visible stereo videos of multiple people.

Key words: Local self-similarity, Mutual information, Multimodal video registration, Dense stereo correspondence, Thermal camera, Visible camera, Visual surveillance

1 Introduction

In the recent years, there has been a growing interest in visual surveillance using thermal-visible imaging system for civilian applications, due to the reduction in the price of infrared sensors. The advantages of jointly using a thermal camera with a visible camera have been discussed comprehensively in [1, 2, 3]. For human monitoring applications in uncontrolled settings, the joint use of these two sensors improves the quality of input data. For example, in a scene where there are shadows on the ground, poor color information under low lighting conditions, or similarity of the human body/clothing with the background, the combined data enables better detection and tracking of people. Moreover, for human activity analysis, the joint use of thermal and visible data enables us to better detect and segment the regions related to the object that people may carry based on their temperature differences compared to the human body.

In the literature, several methods including data fusion algorithms, background subtraction, multi-pedestrian tracking, and classification have been proposed for long-range thermal-visible videos of multiple people [4, 5, 6].

18 However, a fundamental and preliminary task associated with the joint use
19 of thermal-visible data is accurately matching features of a pair of images
20 captured by two different sensors with high differences in imaging character-
21 istics. This task is challenging, especially for close-range scene due to the
22 large scale objects with different detailed patterns in a pair of thermal and
23 visible images. For a pair of close-range videos, it is very difficult to find
24 the correspondence for an entire scene. Therefore, registration is focused on
25 image region of interest (ROI) where for human monitoring applications are
26 human body regions. Matching corresponding regions belonging to a human
27 body in a pair of visible and thermal images is still challenging, because cor-
28 responding pixels have different intensities and have different patterns and
29 textures due to the differences in thermal and visible image characteristics.

30 In few previous related works, MI is the only similarity measure used
31 in dense multimodal stereo matching [7, 8, 9]. Fookes *et al.* proposed a
32 MI-based window matching method that incorporates prior probabilities of
33 the joint probability histogram of all the intensities in the stereo pair in the
34 MI formulation [9]. Their matching method is less sensitive to MI window
35 sizes. However, in their experiment, they only used negative and solarized
36 images that have similar patterns within corresponding ROIs as opposed to
37 thermal and visible images. Egnal has shown that mutual information (MI)
38 is a viable similarity metric for matching disparate thermal and visible im-
39 ages [10]. Chen *et al.* proposed a MI-based registration method for pairs of
40 thermal and visible images with the assumption that each window bound-
41 ing a ROI represents a single human [8]. In their method, occluded people
42 that are merged into one ROI may not be accurately registered since an ROI

43 may contain people within different depth planes. As a solution to improve
44 registration of occluded people in a scene, Krotosky and Trivedi proposed a
45 disparity voting (DV) matching approach [7]. DV is performed by horizon-
46 tally (column by column) sliding small width windows on rectified thermal
47 and visible images, computing MI for pairs of windows, and finally for each
48 column, counting the number of votes associated to each disparity and as-
49 signing one disparity to each column based on a Winner Take All (WTA)
50 approach. Their method can handle occlusion horizontally (two neighboring
51 columns might be assigned to different disparities), but it cannot accurately
52 register people with different height where a shorter person is in front of a
53 taller one (vertical occlusion) since all pixels of a column inside an ROI are
54 assigned to only one disparity.

55 In the abovementioned papers, the correctness and confidence of MI com-
56 pared to other viable similarity metrics is not discussed. Based on our ex-
57 periments, in videos where people have textured clothes, where human ROI
58 segmentation is imperfect (i.e., partial misdetection or false detection), and
59 where there are occlusions, MI is unreliable for matching small width win-
60 dows like the one suggested in [7]. For MI-based stereo matching, choosing
61 the appropriate image window size is not straightforward due to the afore-
62 mentioned difficulties. Also, there is always a trade-off between choosing
63 larger windows for matching evidence, and smaller windows for the precision
64 and details needed for an accurate registration.

65 In this work, we apply local self-similarity (LSS) to the problem of thermal-
66 visible stereo correspondence for close-range human monitoring applications.
67 LSS has been proposed by Shechtman and Irani in [11] and has been pre-

68 viously applied to problems of object categorization, image classification,
69 pedestrian detection, image retrieval by sketching, and object detection [12,
70 13, 14, 11]. To the best of our knowledge, nobody has previously applied
71 LSS as a thermal-visible dense stereo correspondence measure. LSS, simi-
72 larly to MI, computes statistical co-occurrence of pixel intensities. However
73 LSS, unlike MI, is firstly computed and extracted from an individual image
74 as a descriptor and then compared between images. The property of LSS,
75 which makes this measure more interesting for our application, is that the
76 basic unit for measuring internal joint pixel statistics is a small image patch
77 that captures more meaningful image patterns than individual pixels as used
78 in MI computation. This property is useful for matching thermal and visible
79 human ROIs with different direct visual properties such as colors and pat-
80 terns but similar layout/body shape which is an indirect image property. The
81 algorithms presented in this manuscript are based on our previous work [15],
82 but they are further developed with detailed analysis and new evaluations.

83 In section 2, we present a theoretical analysis of LSS and MI as dense
84 multimodal correspondence measures. In section 3, we quantitatively as-
85 sess the reliability and accuracy of MI and LSS as dense stereo similarity
86 measures in various close-range challenging human monitoring scenarios. In
87 section 4, we propose our LSS-based registration that accurately performs
88 for multiple people and occlusions. Finally, in section 5, we qualitatively
89 and quantitatively compare our LSS-based stereo registration method and a
90 recent state-of-the-art MI-based stereo registration method for human mon-
91 itoring applications.

92 **2. Theoretical analyses of MI and LSS for thermal-visible human**
93 **ROI correspondence**

94 *2.1. Mutual information*

95 Mutual information (MI) is the classic dense similarity measure for mul-
96 timodal stereo registration. The MI between two image windows L and R is
97 defined as

$$MI(L, R) = \sum_l \sum_r P(l, r) \log \frac{P(l, r)}{P(l)P(r)}, \quad (1)$$

98 where $P(l, r)$, is the joint probability mass function and $P(l)$ and $P(r)$ are the
99 marginal probability mass functions. $P(l, r)$ is a two-dimensional histogram
100 $g(l, r)$ normalized by the total sum of the histogram. $g(l, r)$ is computed as
101 for each point, the quantized intensity levels l and r from the left and right
102 matching windows (L and R) increment $g(l, r)$ by one. The marginal proba-
103 bilities $P(l)$ and $P(r)$ are obtained by summing $P(l, r)$ over the grayscale or
104 thermal intensities.

105 The unit of measure for MI is the pixel, which forces existing underlying
106 visual properties (i.e., pixel colors, intensities, edges, or gradients) in thermal
107 and visible images to be identical for a contribution in MI computation. In
108 our application, MI computes the statistical co-occurrence of pixel-wise mea-
109 sures, such as patterns related to human body regions on pairs of thermal
110 and visible images. Based on our experiments, MI is unreliable for matching
111 differently textured corresponding human body ROIs and partially misde-
112 tected or falsely detected human body ROIs caused by erroneous foreground
113 segmentation in thermal and visible images. MI only performs well when
114 the joint probability histogram is sufficiently populated inside MI windows.

115 Choosing the appropriate window size is not straightforward due to the afore-
 116 mentioned difficulties. In fact, because of imperfect data, the appropriate size
 117 might not be available for the computation (e.g. region fragmentation, small
 118 objects).

119 2.2. Local self-similarity

120 Local self-similarity (LSS) is a descriptor that capture locally internal
 121 geometric layout of self-similarities (i.e., edges) within an image region (i.e.,
 122 human body ROI) while accounting for small local affine deformation. Ini-
 123 tially, this descriptor has been proposed by Sechtman and Irani [11]. LSS
 124 describes statistical co-occurrence of small image patch (e.g. 5×5 pixels) in a
 125 larger surrounding image region (e.g. 40×40 pixels). First, a correlation sur-
 126 face is computed by a sum of the square differences (SSD) between a small
 127 patch centered at pixel p and all possible patches in a larger surrounding
 128 image region. SSD is normalized by the maximum value of the small image
 129 patch intensity variance and noise (a constant that corresponds to acceptable
 130 photometric variations in color or illumination). It is defined as

$$S_p(x, y) = \exp\left(-\frac{SSD_p(x, y)}{\max(var_{noise}, var_{patch})}\right). \quad (2)$$

131 Then, the correlation surface is transformed into a log-polar representation
 132 partitioned into e.g. 80 bins (20 angles and 4 radial intervals). The LSS
 133 descriptor is defined by selecting the maximal value of each bin that results
 134 in a descriptor with 80 entries. A LSS descriptor is firstly computed for a
 135 ROI within an image then it can be compared with other LSS descriptors in
 136 a second image using a measure such as $L1$ distance.

137 Previously, Shechtman and Irani have shown that for matching image
138 regions with similar shape/layout but with different direct visual properties
139 such as color and edges, LSS is a more reliable similarity metric compared to
140 other local image descriptors and match measures such as MI [11]. Shecht-
141 man and Irani have also shown that LSS is applicable for image retrieval
142 by sketching [11]. In their work, the template image is a sketch of human
143 body representing a body pose. The template is used to detect similar hu-
144 man body poses in several images with different color and textures. For this
145 application, the advantage of LSS is that its unit measure, which is a small
146 image patch, contains more meaningful patterns compared to pixel as used
147 for MI computations. As it is described in Shechtman and Irani’s work [11],
148 this property makes LSS a suitable measure for matching image regions with
149 different direct visual properties such as color, edges, or textures, as long as
150 they have similar spatial layouts.

151 For matching thermal and visible human ROIs, we believe that the same
152 LSS property as used for image retrieval by sketching in [11] is applicable. In
153 fact, in thermal and visible images, the edges and pixel intensities within cor-
154 responding human body ROIs are not identical as it is required for matching
155 using MI, but the human body layout is a common indirect visual property
156 between thermal and visible corresponding human ROIs.

157 Registration accuracy is an important factor for human ROIs registration
158 in a scene with multiple people and occlusions. In order to apply LSS as
159 a similarity metric, it is required that this descriptor captures local image
160 ROI layout within a small surrounding region while considering the required
161 image details. Therefore, we experimentally found out that for a close-range

162 video, a patch of size 3×3 as a unit measurement within surrounding region
 163 of 20×20 pixels is sufficient to capture meaningful local image patterns of
 164 human body shape that are mostly belonging to edges. For an LSS-based
 165 window matching similar to MI-based matching, the overall geometric layout
 166 within a matching window is captured by a set of LSS descriptors respecting
 167 the relative geometric positions of the descriptors. Although the window
 168 sizes have a width smaller than the width of a complete human body and a
 169 height the same as the human body height; the set of descriptors within a
 170 window still captures partially body geometric layout. In fact, matching is
 171 performed by comparing two sets of descriptors belonging to two matching
 172 windows in thermal and visible images. A good match corresponds to two
 173 windows where the descriptors are similar both in values and their relative
 174 geometric positions.

175 In order to perform a better matching, we discard the non-informative
 176 descriptors from each set of descriptors. Non-informative descriptors are
 177 the ones that do not contain any self-similarities (i.e., the center of a small
 178 image patch is salient) and the ones that contain high self-similarities (i.e., a
 179 homogenous region with a uniform texture/color). A descriptor is salient, if
 180 all its bins' values are smaller than a threshold. The homogeneity (which also
 181 cause a non-informative descriptor) is detected using the sparseness measure
 182 of [16]. The sparseness measure is defined as

$$\text{sparseness}(X) = \frac{\sqrt{n} - (\sum |x_i|) / \sqrt{\sum x_i^2}}{\sqrt{n} - 1} \quad (3)$$

183 where n is the dimensionality of descriptor x (in our method 80). This func-
 184 tion evaluates to unity if and only if x contains only a single non-zero compo-
 185 nent, and takes a value of zero if and only if all components are equal. Dis-

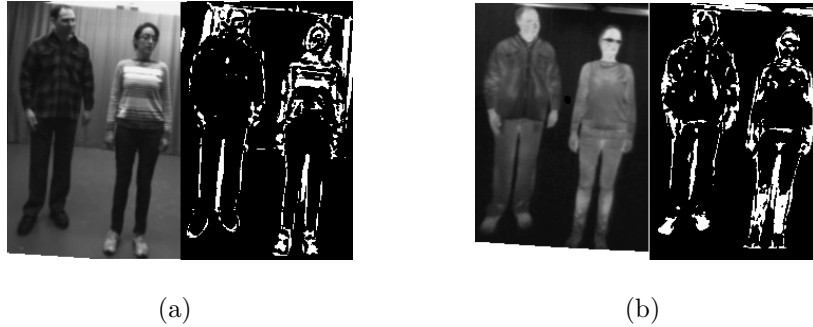


Figure 1: Informative LSS descriptors. (a) Visible and informative LSS descriptors images
 (b) Thermal and informative LSS descriptors images.

186 carding non-informative descriptors is like an implicit segmentation or edge
 187 detection, which for window matching, increases the discriminative power
 188 of the LSS measure and avoids ambiguous matching. It is important to
 189 note that the remaining informative descriptors still form a denser collection
 190 compared to sparse interest points. Fig. 1 shows pixels having informative
 191 descriptors (white pixels) for a pair of thermal and visible images. Fig. 1
 192 is an extreme example of discarding non-informative pixels to highlight by
 193 filtering out those pixels, the common visual property between thermal and
 194 visible human ROIs (i.e., human body shape) are obtainable without any
 195 explicit edge detection or segmentation.

196 2.3. *Introductory examples of human ROI matching*

197 In order to illustrate the difficulties of thermal-visible human ROI match-
 198 ing and the advantages of LSS compared to MI, we present three introductory
 199 examples of human ROI matching using a simple sliding window matching
 200 approach and various window sizes. In these examples, matching is per-
 201 formed by computing the similarity distances of a fixed window on an image

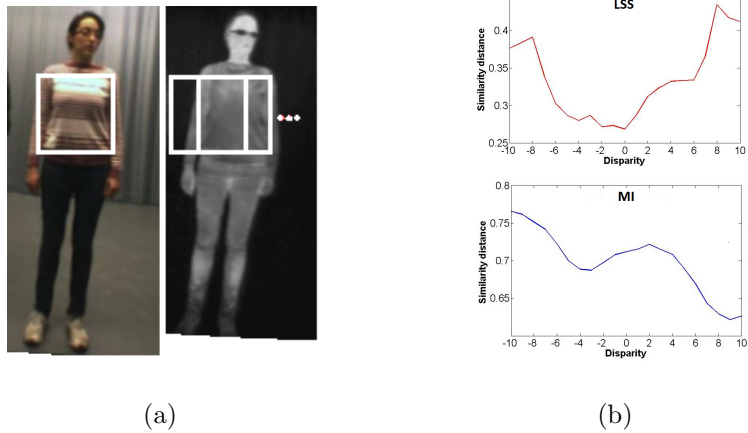


Figure 2: Matching corresponding textured and uniform regions in visible and thermal pair of images. (a) Aligned visible and thermal images and (b) Similarity distances of LSS and MI for disparity interval of $[-10,10]$.

202 ROI of the visible image with a sliding window on the thermal image within
 203 a disparity range of $[-10, 10]$, and then choosing the disparity that minimizes
 204 the similarity distance. In order to simplify the search to 1D, the two images
 205 were rectified, and then manually aligned so that a disparity of 0 corresponds
 206 to a ground-truth alignment (more details about multimodal camera calibra-
 207 tion in section 3.1). We defined the LSS-based similarity distance between
 208 two windows L and R by the sum of the $L1$ distances of informative de-
 209 scriptors within those two windows, and the MI-based similarity distance as
 210 $1 - MI(L, R)$. Fig. 2 shows an example of matching a textured region in the
 211 visible image with a corresponding uniform region in the thermal image. Fig.
 212 2 (b) shows the similarity distance results for both MI and LSS over a preset
 213 disparity range. For LSS, the similarity distance is correctly minimized at
 214 disparity 0. However for MI, the similarity distance is minimized incorrectly.

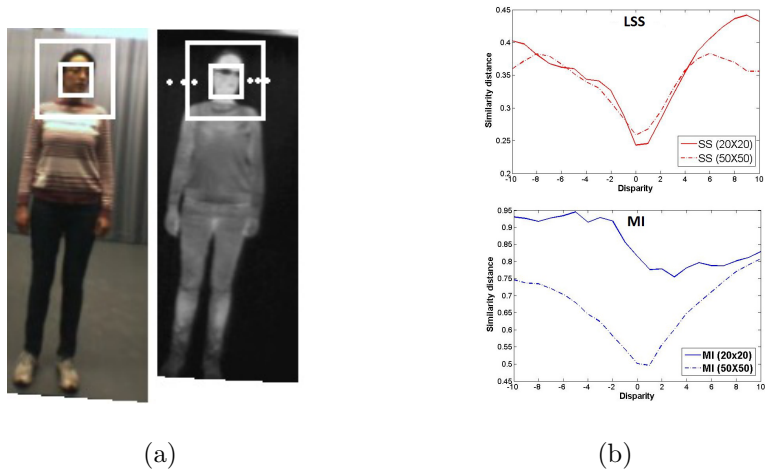


Figure 3: Matching corresponding regions of visible and thermal within image windows of size 20×20 and 50×50 pixels. (a) Aligned visible and thermal images, (b) Similarity distances of LSS and MI for disparity interval of $[-10,10]$.

215 This illustrate that MI is not a robust similarity metric for matching a tex-
 216 tured region and a uniform region when there are not many similar patterns.
 217 Fig. 3 shows an example of matching windows of sizes 20×20 and 50×50
 218 pixels on a head region. Fig. 3 (b) shows that MI is not a robust mea-
 219 sure for matching 20×20 thermal-visible windows. However, using larger
 220 window of size 50×50 pixels containing more similar patterns and more
 221 similar spatial layout, MI-based similarity distance is correctly minimized
 222 at disparity 0. For this example, LSS-based similarity distance is correctly
 223 minimized at disparity 0 for both matching window sizes which illustrate the
 224 robustness of this measure for matching even with small window sizes. Fig.
 225 4 shows an example of matching thermal-visible windows on regions with
 226 dramatic partial ROI misdetection using matching window sizes of 20×170
 227 and 60×170 pixels. In the visible image, due to the color similarity of the

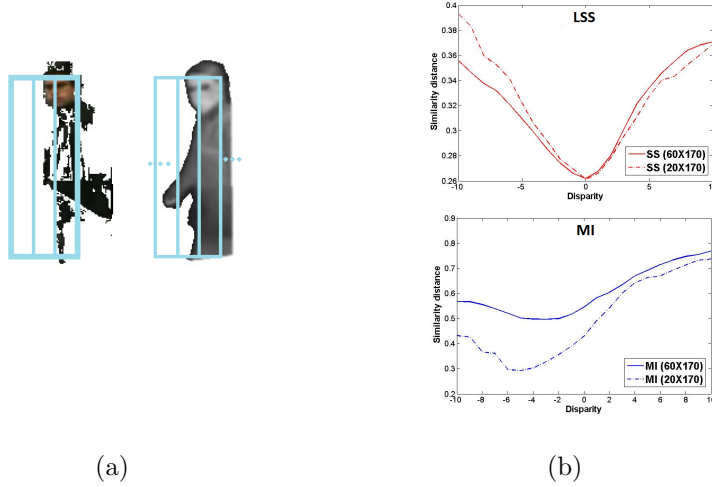


Figure 4: Matching corresponding foreground pixels within 20×170 and 60×170 pixels windows in visible and thermal pair of images (a) Aligned visible and thermal images, (b) Similarity distances of LSS and MI for disparity interval of $[-10,10]$.

228 ROI and the background, some parts of the body region are not detected.
 229 Fig. 4 (b) shows that MI fails to find the correct disparity offset with both
 230 window sizes. However, LSS find the correct disparity which illustrates the
 231 robustness of this measure for partial ROI misdetection.

232 3. Quantitative analyses of MI and LSS for thermal-visible human 233 ROI correspondence

234 In our previous work, we have evaluated the performance of several local
 235 image descriptors and similarity measures for thermal-visible human ROI
 236 registration [17]. In this work, we only focus on the comparison between MI
 237 and LSS with new evaluation criteria and extensive experiments.

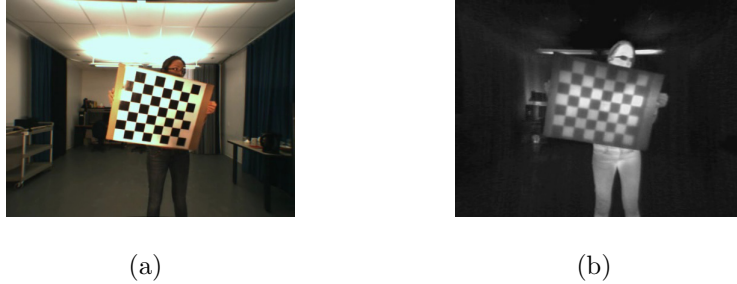


Figure 5: Calibrating images: (a) Visible image and (b) Thermal image.

238 3.1. Video Acquisition and Calibration

239 We used synchronized visible-thermal videos of a $5m \times 5m$ room at a fixed
 240 temperature of $24\text{ }^{\circ}C$ captured by stationary thermal and visible cameras
 241 with a 12 cm baseline. We used sets of video frames of a relatively close range
 242 scene where different people with different poses and clothing are walking at
 243 different depths (between 2-5 meters) from the camera baseline. In order to
 244 simplify the stereo matching to a 1D search, we first calibrated the thermal
 245 and visible cameras, and then rectified the images using the intrinsic and
 246 extrinsic calibration parameters. We used the standard technique available in
 247 the camera calibration toolbox of MATLAB ([18]). For calibration, we placed
 248 a checkboard pattern in front of the cameras. Since in the thermal images,
 249 the checkboard pattern is not visible at room temperature; we illuminated
 250 the scene using high intensity halogen bulbs placed behind the two cameras.
 251 This way, the dark squares absorb more energy and visually appear brighter
 252 than the white squares. Fig. 5 shows an example of our calibration images.

253 3.2. *Experimental setup*

254 Our experimental setup is designed to study the efficiency of MI and LSS
255 as similarity measures for thermal-visible human ROI registration using real-
256 istic videos of multiple people in a close range scene. In our experiment, the
257 ROIs within corresponding windows in a pair of thermal and visible images
258 might be differently textured or one textured and the other uniform. Win-
259 dows are centered at randomly picked points that are located inside visible
260 human ROIs. We used a sliding window matching (see section 3.3) to find
261 the corresponding image window on the thermal image. The matching pro-
262 cess was repeated using three rectangular window sizes of 10×130 (small),
263 20×130 (medium), and 40×130 (large) pixels. The heights of the windows
264 are chosen as the maximum possible height of a person in our experimental
265 videos. The randomly picked points were located either on textured or tex-
266 tureless visible human ROI for relatively near targets (between 2 to 3 meters
267 from the camera) or far targets (between 4 to 5 meters). Note that for close-
268 range scene monitoring, the scale of targets considerably changes by walking
269 one meter further away or toward the camera. Fig. 6 shows an example of
270 randomly picked matching window. Our experiment is carried out using 300
271 matching windows (100 points using three window sizes).

272 3.3. *Sliding window matching*

273 For each thermal and visible pair of images, a window centered at a point
274 on the human ROI at column j on the visible image is defined ($W_{l,j}$). Then, a
275 1D window matching search is done on the thermal image in order to find the
276 corresponding window $W_{r,j+d}$ which minimizes a similarity distance SD . d
277 is a disparity offset belonging to disparity interval set D . In our experiment,



Figure 6: Thermal-visible 1-D sliding window matching.

278 the size of D is the same size as the image width. Fig. 6 illustrates the
 279 sliding window matching.

280 For LSS, the descriptor computation and the matching are done in two
 281 separate processes, for each pair of image windows $W_{l,j}$ and $W_{r,j+d}$ centered
 282 at column j on the visible image and column $j + d$ on the thermal image.
 283 A normalized similarity distance $SD_{j,d}$, which is the sum of $L1$ distance
 284 of the corresponding pixels $p_l \in W_{l,j}$ and $p_r \in W_{r,j+d}$ having informative
 285 descriptors, is computed as

$$SD_{j,d} = \frac{\sum_{p_l, p_r} L1_{l,r}(p_l, p_r)}{N}, \quad (4)$$

286 where N is the number of corresponding pixels p_l and p_r contributing in the
 287 similarity distance computation and d is the disparity offset. This number is
 288 also proportional to the number of informative pixels inside an image ROI.
 289 The typical value of N for window size of 40×130 varies in the range of 200
 290 to 1000 pixels and it is maximum when edges and boundaries inside matching
 291 windows are correctly overlapped. $L1_{l,r}$ is computed as

$$L1_{l,r}(p_l, p_r) = \sum_{k=1}^{80} |d_{p_l}(k) - d_{p_r}(k)| \quad (5)$$

292 where 80 is the number of local self-similarity descriptor bins.

293 For MI, SD is defined as

$$SD_{j,d} = 1 - MI(W_{l,j}, W_{r,j+d}), \quad (6)$$

294 where MI is the mutual information defined in equation 1. And finally the

295 best disparity associated to best matching windows is computed by

$$d_{min} = \underset{d}{\operatorname{argmin}} (SD_{j,d}), d \in D. \quad (7)$$

296 3.4. Evaluation Criteria

297 In our evaluation, we assess the precision-recall and power of discrimina-
298 tion of MI and LSS as explained in the following sections.

299 3.4.1. Precision and recall

300 We used a criterion similar to the one used in [19]. Precision and recall
301 are defined as follows:

$$\textit{precision} = \frac{\#correctmatches}{\#matchesretrieved} \quad (8)$$

302

$$\textit{recall} = \frac{\#correctmatches}{\#totalcorrespondences} \quad (9)$$

303 In our experiment, *correctmatches* is the number of matches with a dis-
304 parity error smaller than 3 pixels with respect to ground-truth and with SD
305 (equation 4 or 6) smaller than a threshold t (t varies between minimum pos-
306 sible values where *matchesretrieved* become one and maximum value where
307 *matchesretrieved* become all the matched windows *totalcorrespondence*).
308 *totalcorrespondence* is a fixed value that corresponds to the number of tested
309 windows (i.e., 100 windows of each size). *matchesretrieved* is the number

310 of matches with a SD below threshold t . $matchesretrieved$ varies from 1 to
311 $totalcorrespondences$. In a precision and recall curve, a feature with high
312 recall value and low precision value means that many correct matches as well
313 as many false matches are retrieved. On the other hand, high precision value
314 and low recall value means that most matches are correct but many others
315 have been missed.

316 3.4.2. Power of discrimination

317 To assess the reliability of a similarity metric, not only its precision is im-
318 portant but also how that similarity metric possesses isolation characteristic
319 (power of discrimination) is important as well.

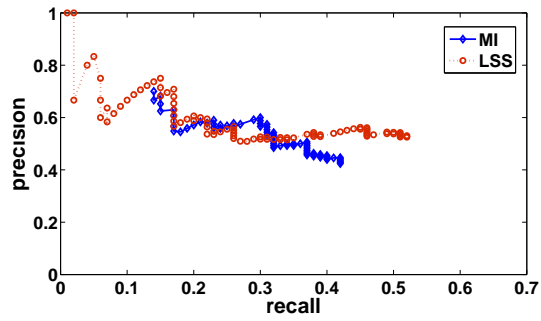
320 A similarity metric possesses a high power of discrimination, if its correct
321 matches are located on isolated minimums over D (disparity range) and
322 SD (equation 4) curve (that is, having SD value much smaller than its
323 neighbors). In order to evaluate the isolation characteristic of MI and LSS,
324 for their correct matches, we study the shape of SD computed along the
325 disparity range $D = [q - 20 : q + 20]$, where q is the position of the global
326 minimum (best match). We applied the same measure as in [20]. In order
327 to evaluate the isolation of the global minimum, the SD values computed by
328 the sliding window matching (section 3.3) are first sorted increasingly and
329 are transformed to the interval $[0, 1]$ named SD' . Second, N is the number of
330 values in SD' that are less than a pre-computed small threshold α , ignoring
331 the global minimum. α has the same value for evaluating all descriptors
332 and measures. Third, a quality measure s (the s value) is computed by
333 dividing N by the size of the disparity range. So $s = 0$ corresponds to
334 the most isolated minimum (best performance), and $s = 1$ corresponds to

335 the least isolated minimum (flat/constant SD versus d curve). Finally, for
336 each correspondence measure, a graph of Accumulated Frequencies (AF)
337 of the s values of all matches is computed (In fact AF is the distribution
338 of s values belonging to correct matches). Therefore, the correspondence
339 measure for which AF reaches a higher value at a smaller s value is the more
340 discriminative.

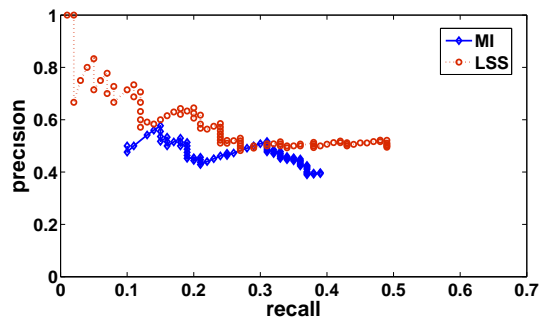
341 3.5. Results

342 First, we present the evaluation of MI and LSS using the precision and
343 recall criterion as explained in section 3.4.1. Fig. 7 shows the precision-recall
344 curves of MI and LSS for small, medium, and large window sizes as described
345 in section 3.2. Overall, for all the three matching window sizes, LSS achieves
346 higher values of recall and precision compared to MI. The largest size window
347 achieves better precision than medium and small ones for both MI and LSS.
348 However, MI is totally inefficient for small window sizes. This result shows
349 the robustness of MI is closely related to the sizes of MI windows, which
350 are required to be large enough to sufficiently populate the joint probability
351 histogram. On the other hand, the precision of LSS is more consistent for
352 three window sizes. For this experiment, we used a simple sliding window
353 matching that ignores the occlusions. Using a more appropriate correspon-
354 dence algorithm that we propose in section 4, will result in higher matching
355 precisions.

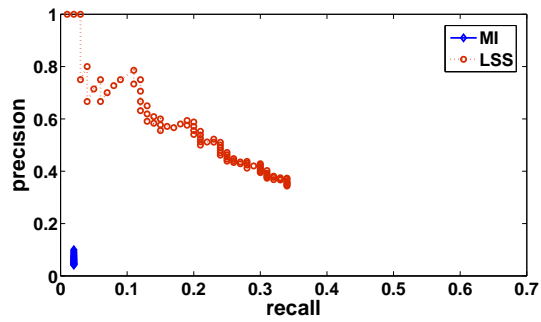
356 Fig. 8 shows the accumulated frequency distribution of s (details in
357 section 3.4.2) obtained for MI and LSS using three window sizes as described
358 in section 3.2. It can be seen, for LSS compared to MI, AF starts with a
359 higher value and reaches to the higher value for a smaller S value. This shows



(a)

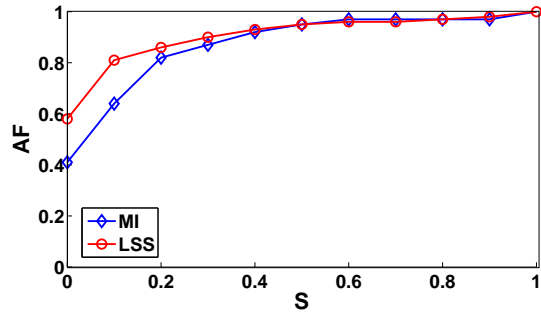


(b)

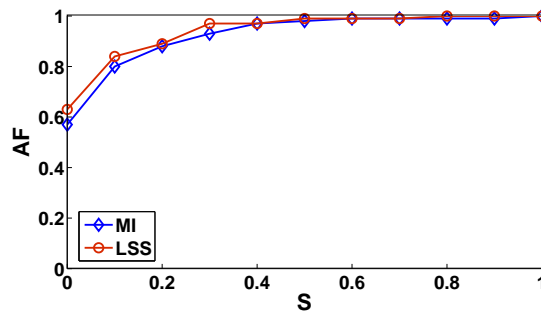


(c)

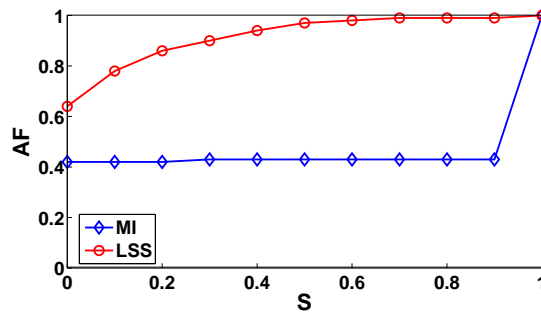
Figure 7: Precision- recall curves: (a) large window (40×130) (b) medium window (20×130) (c) small window (10×130).



(a)



(b)



(c)

Figure 8: Accumulated frequencies versus S value: (a) large window (40×130) (b) medium window (20×130) (c) small window (10×130).

360 LSS possesses a better isolation characteristics compared to MI.

361 Overall, the results show that comparing to MI, LSS is a more reliable
362 similarity metric for matching differently textured human ROIs in thermal
363 and visible images and it is less restricted by size of matching windows.

364 4. LSS-based multimodal ROI registration

365 In this section, we describe our novel multimodal ROI registration method
366 using LSS. For a pair of thermal and visible video frames, our goal is to
367 register the ROIs belonging to moving people in a scene in which they may
368 be temporary stationary for a few frames. Our method addresses registration
369 of multiple people merged into one ROI with different levels of occlusion and
370 with partially erroneous foreground segmentation for realistic thermal-visible
371 videos of a close range scene. We assume that each person at each instant lies
372 approximately within one depth plane in the scene. Therefore, we propose
373 that a natural way for estimating depth planes related to multiple moving
374 people is by applying motion segmentation on foreground pixels with the
375 assumption that each motion segment belongs to one person in the scene,
376 but more than one motion segment may belong to a person.

377 We define the multimodal image registration as multiple labeling sub-
378 problems. Then, we use the disparity voting matching approach to register
379 each individual motion segment rather than a whole foreground blob. Let
380 MS be the set of motion segments belonging to moving people in the scene,
381 and D be a set of labels corresponding to disparities. Our registration method
382 assigns a label $d_k \in D$ in the range between d_{min} to d_{max} to each pixel
383 of a motion segment $ms_i \in MS$. Thus, our registration method has two

384 main parts: 1) motion segmentation that divides the registration problem
 385 as multiple labeling sub-problems and 2) disparity assignment which assigns
 386 disparity to each segment. The two parts of our method are described in the
 387 subsequent sections.

388 *4.1. Motion segmentation*

389 Our motion segmentation has three steps. Firstly, we extract foreground
 390 pixels using the background subtraction method proposed in [21]. Any back-
 391 ground subtraction method with a reasonable amount of error is applicable.
 392 Secondly, we compute the motion vector field for foreground pixels using an
 393 optical flow method based on block-matching [22]. To speed up the process,
 394 the optical flow is only computed for regions inside the bounding boxes of
 395 the union of the foreground masks of two consecutive frames $t - 1$ and t ,
 396 instead of the whole image. Thirdly, we apply the mean-shift segmentation
 397 method proposed in [23] for segmenting the motion vector fields computed
 398 in the previous step and computing a mean velocity vector for the computed
 399 segments. Mean-shift segmentation is applied on $(2+2)$ feature point dimen-
 400 sions, where two dimensions are related to spatial dimensions (horizontal and
 401 vertical directions) and the two others are related to the two motion vector
 402 components in x and y directions. Applying motion segmentation on ROIs
 403 results in a set of motion segments S defined as

$$SM = \{sm_1, \dots, sm_m\}. \quad (10)$$

404 An average mean velocity vector \hat{m}_i is associated to each sm_i using

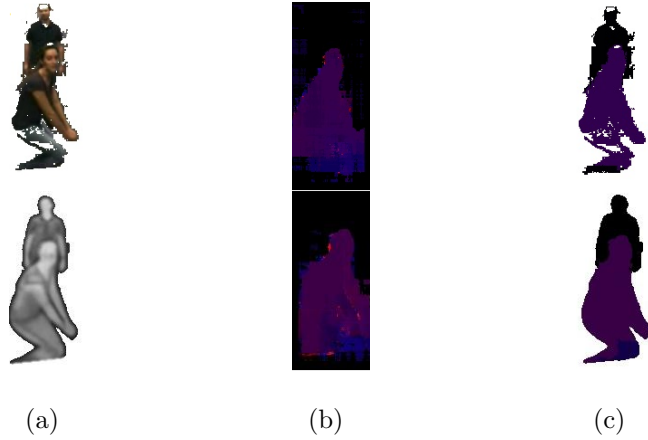


Figure 9: First row: visible image and second row: thermal image. (a) Foreground image, (b) motion field vectors, and (c) motion segmentation (depth segments).

$$\hat{m}_i = \frac{\sum_{p \in sm_i} m(p)}{|sm_i|}, \quad (11)$$

405 where $m(p)$ is the motion vector of pixel p . Fig. 9 shows the motion
 406 segmentation results of two occluding people, where one of them is temporary
 407 stationary. Motion vectors are visualized by a mapping to HSV color space.
 408 Applying motion segmentation on foreground pixels enables us to determine
 409 also a depth segment associated to temporary stationary person for which
 410 its mean velocity vector is zero. Since in most indoor videos, the motion
 411 segmentation of thermal images are more accurate compared to visible images
 412 due to less partial ROI misdetection error, we perform motion segmentation
 413 for thermal images and we register the thermal motion segments on visible
 414 foreground images. However, it could also be done the opposite way.

4.2. Disparity assignment

At this step, we assign disparity to each motion segment individually. We use a disparity voting matching approach similar to the one that was previously proposed by Krotosky and Trivedi [7]. DV matching assigns one single disparity to all the pixels of a column of matching regions. However, different disparities can be assigned to two neighboring columns. Krotosky and Trivedi DV method uses MI as similarity metric and is performed on whole foreground blobs. Their method is able to resolve the horizontal part of an occlusion, but fails to assign correct disparity for the vertical part of an occlusion (in this case, the pixels of a column for a region associated to vertically occluded people should be assigned to a different disparity) (see fig. 11). To solve this problem, we propose performing DV on each motion segment separately. Moreover, based on our previous experiments, we use the informative LSS descriptors as similarity measure.

4.2.1. LSS-based DV algorithm

For each $sm_i \in S$, we build a disparity voting matrix of DV_i of size $(N, d_{max} - d_1 + 1)$ where N is the number of pixels of sm_i and $[d_1 - d_{max}]$ is a preset disparity range. This procedure is performed by shifting column by column $W_{l,j}$ on the reference image for all the columns $j \in s_i$, then doing window matching, the same as we previously described in section 3.3. Then, for each d_{min} computed by window matching, a vote is added to $DV_i(p_l, d_{min})$ for all $p_l \in (W_{l,j} \cap s_i)$. Since the width of windows are m pixels wide, we have m votes for each pixel belonging to s_i . Finally, the disparity map DM_i is computed as,

$$DM_i(p_l) = \operatorname{argmax}_d (D_i(p_l, d)), \quad (12)$$

439 5. Experimental validation and discussion

440 We have assessed our registration method with over 5000 video frames
 441 of up to 5 people with different clothing, various poses, distances to cam-
 442 eras, and with different level of occlusions. In these experiments, we used
 443 the same experimental setup as described previously in section 3.1. The first
 444 test video was captured during summer with people having lighter clothes
 445 (light clothes results in less heat patterns on the body in infrared) and with
 446 a fair amount of textures inside human ROIs in thermal and visible im-
 447 ages. The background subtraction errors were mostly misdetection errors.
 448 Our other two test videos were captured during winter with people wearing
 449 winter clothes (thick clothes results in more heat patterns on the body in
 450 infrared), which causes patterns inside human body ROIs. The background
 451 subtraction results in our winter videos include both misdetection errors and
 452 falsely detected region as foreground. The disparity range was between 5 to
 453 50 pixels.

454 Fig. 10 illustrates successful registrations with our method in one of our
 455 winter videos for three frames of people in different levels of occlusions.

456 5.1. Comparison of DV correspondence and our correspondence algorithm

457 In order to demonstrate the accuracy improvement of our method com-
 458 pared to a state-of-the-art disparity voting algorithm (DV) [7] in handling
 459 occlusions, we quantitatively compared our disparity results using motion



Figure 10: Registration results of foregrounds using imperfect background subtraction with false positive and false negative errors.

460 segmentation and the results of DV using for both LSS as similarity mea-
 461 sure. We generated ground-truth disparities by manually segmenting and
 462 registering regions of foreground for each frame. Fig. 11 illustrates the
 463 comparison with ground-truth. Results in the first and second rows illus-
 464 trate examples where two people in two different depths in the scene are
 465 in occlusion. LSS+DV method fails to assign correct different disparities
 466 to the columns containing pixels related to more than one individual since
 467 based on a WTA approach, a single disparity is assigned to all the pixels of
 468 each column. However, LSS+MS+DV succeed in assigning accurately dif-
 469 ferent disparities to the two human body ROIs since the DV was applied to
 470 each motion segment individually. Accordingly, in fig. 11 (d), for the first
 471 and second rows, the sum of disparity errors of the columns corresponding
 472 to two occluded people is much higher for LSS+DV method compared to
 473 LSS+MS+DV method.

474 Indeed, to register merged objects in a single region, DV makes no as-
 475 sumptions about the assignment of pixels to individual objects and assigns

476 a single disparity to each column inside an ROI based on a maximization of
477 the number of votes. In their matching approach [7], if a column of pixels
478 belongs to different objects at different depth in the scene, the vote only
479 goes for one of them based on WTA approach. However, in our registration
480 method, motion segmentation gives a reasonable estimate of moving regions
481 belonging to people in the scene, and applying the DV matching on each mo-
482 tion segment gives more accurate results since it is less probable that pixels
483 in one column belongs to more than one object. Therefore, in the worst case,
484 even with erroneous motion segmentation, our method will have at minimum
485 the same accuracy as the DV algorithm.

486 Fig. 11, last row, illustrates the example of multiple occluding people. Al-
487 though LSS+MS+DV registration results are not perfect because few small
488 motion segments resulting from over segmentation were not matched cor-
489 rectly still the results are more accurate than for LSS+DV. Accordingly, in
490 Fig. 11 (d), last row, the sums of disparity error for columns related to verti-
491 cal occlusion is higher for LSS+DV compared to LSS+MS+DV. However, it
492 is noticeable that in some columns, LSS+MS+DV has slightly higher errors
493 caused by small motion segments misalignment.

494 Fig. 12 illustrates other registration results with LSS+MS+DV and
495 LSS+DV. It is observable, that for LSS+DV method, the object misalign-
496 ments happen where there are vertical occlusions while our method performs
497 accurately in such a case.

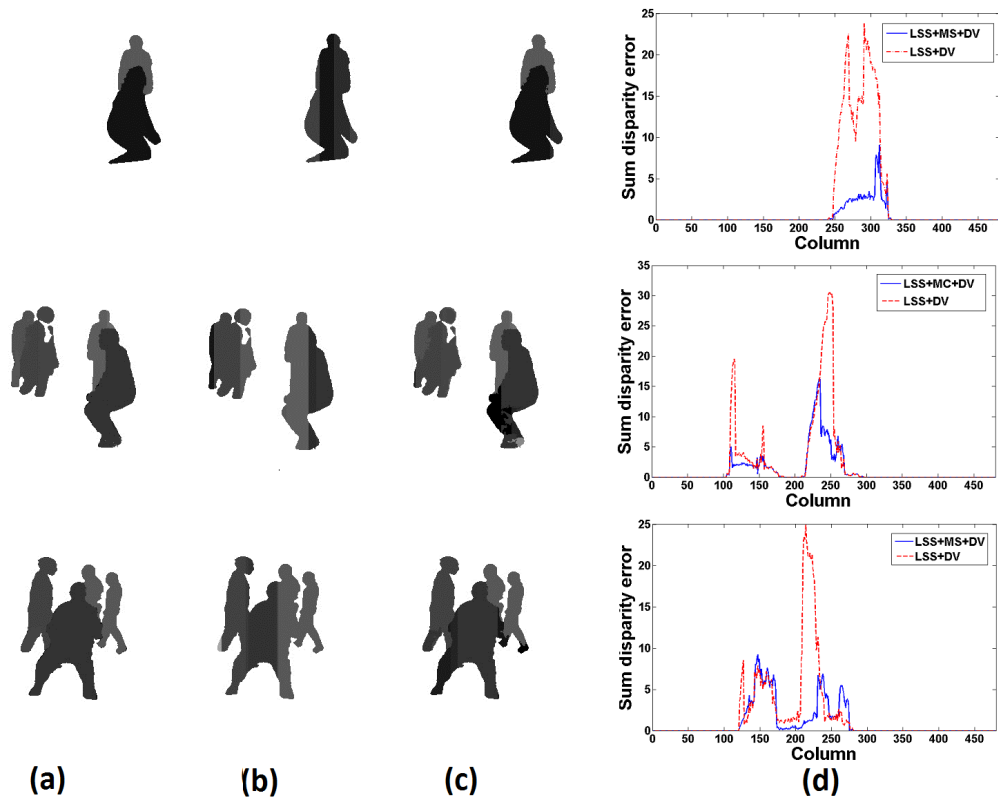


Figure 11: Comparison of LSS-based DV method and our proposed disparity assignment method (a) ground-truth disparity, (b) disparity estimation of DV matching using LSS as similarity measure (LSS+DV), (c) disparity estimation of our proposed method (LSS+MS+DV), and (d) Sum disparity errors over each column of pixels.

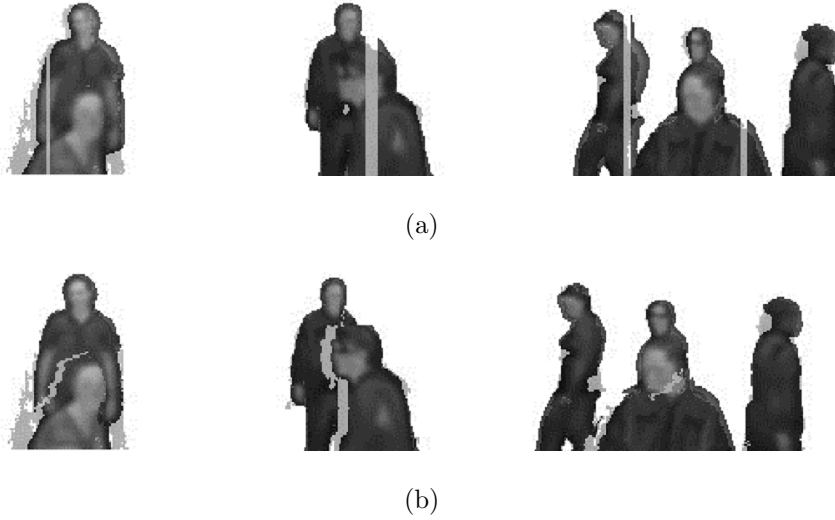
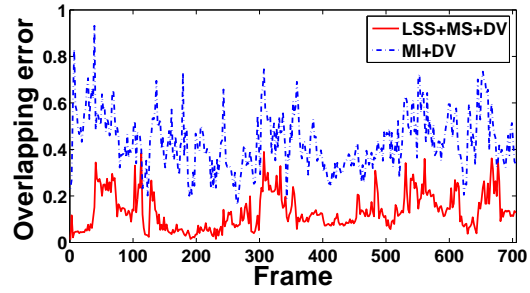


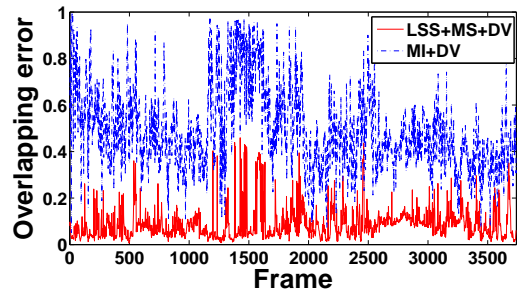
Figure 12: Comparison of LSS+DV and LSS+MS+DV detailed registration: (a) LSS+DV registration and (b) LSS+MS+DV registration.

498 *5.2. Comparison of our LSS-based registration with the state-of-the-art MI-*
 499 *based registration*

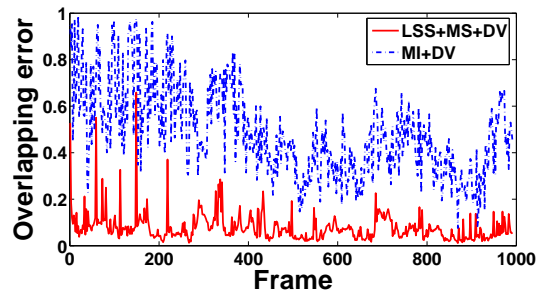
500 In order to demonstrate the improvement of our LSS-based registration
 501 method LSS+MS+DV compared to the state-of-the-art MI-based registra-
 502 tion method, MI+DV, proposed by Krotosky and Trivedi [7], we qualitatively
 503 and quantitatively compared the two methods. Fig. 14 illustrates four exam-
 504 ples of the disparity computation and the image registration results obtained
 505 using the two methods for our summer video. Note that our results are more
 506 accurate, especially for occlusions. Fig. 15 illustrates four examples for win-
 507 ter video 1. Note that MI+DV results are significantly poorer. These results
 508 demonstrate that for videos where there are falsely detected region as fore-
 509 ground and high differences of patterns inside human body ROIs, MI is not
 510 a reliable similarity measure. In contrast, LSS performs very well, except for



(a)



(b)



(c)

Figure 13: Overlapping error: (a) Summer video (702 frames), (b) Winter video 1 (3740 frames), and (c) Winter video 2 (992 frames)

511 few misalignments which occur for very small motion segments.

512 For a quantitative evaluation of the two registration methods, we defined
513 an overlapping error that gives a quantitative estimate of the registration
514 accuracy. The overlapping error is defined as,

$$E = 1 - \frac{N_{v \cap t}}{N_t}, \quad (13)$$

515 where $N_{v \cap t}$ is the number of overlapping aligned thermal foreground pix-
516 els on visible foreground pixels and N_t is the number of thermal foreground
517 pixels. The best performance with zero overlapping error is when all the
518 thermal pixels on the reference image have corresponding visible pixels on
519 the second image. Note that our registration results are aligned thermal on
520 visible images. This evaluation measure includes the background subtraction
521 errors and also ignores misaligned thermal pixels which have falsely matched
522 visible foreground pixels. However, since for both methods the background
523 subtraction errors are included in the overlapping error, the differences be-
524 tween the two methods errors are still a good indicator for comparing overall
525 registration accuracies for a large numbers of frames. Fig 13 illustrates the
526 overlapping error using our LSS+MS+DV and MI+DV [7] methods for sum-
527 mer and winter videos. Based on table 1, the differences of mean overlapping
528 error for the two methods over all frames (fig 13 (a)) are 0.30 for the sum-
529 mer video (fig 13 (a)), and 0.40 and 0.41 for the winter videos, (fig 13 (b)
530 and (c), respectively). Also, table 1 shows the standard deviation of over-
531 lapping. The fluctuation of overlapping error for LSS+MS+DV method is
532 much less than for MI+DV [7] method, especially for winter videos because
533 of larger difference in textures on the objects. These results demonstrate

534 that our method performs more accurately and more consistently compared
 535 to MI+DV [7] method, especially for winter videos, in accordance with our
 536 qualitative results and previous discussions.

Table 1: Overlapping error (OE) for disparity voting (MI) and our proposed algorithm (LSS) with multiple people in the scene: frames with occlusion. SV: summer video, WV1 and WV2: winter videos, NO: number of objects, NF: number of frames, SM: similarity metric, and % OE (Ave - Std): average and standard deviation of overlapping error.

Video	NO	NF	SM	OE (Ave - Std)
SV	4	702	LSS	0.13 - 0.07
			MI	0.43 - 0.12
WV1	4	3740	LSS	0.09 - 0.06
			MI	0.49 - 0.17
WV2	5	992	LSS	0.07 - 0.07
			MI	0.48 - 0.19

537 6. Conclusion

538 In this paper, we applied LSS as a multimodal dense stereo correspon-
 539 dence measure and shown its advantages compared to MI, the most com-
 540 monly used multimodal stereo correspondence measure in the state-of-the-
 541 art for human monitoring applications. We also proposed an LSS-based
 542 registration method, which addresses the accurate registration of regions as-
 543 sociated to occluded people in different depths in the scene. In our results,
 544 we have shown the improvement of our registration method over the DV
 545 method proposed by [7]. Moreover, we have shown that our method signifi-
 546 cantly outperforms the state-of-the-art MI-based registration method in [7].

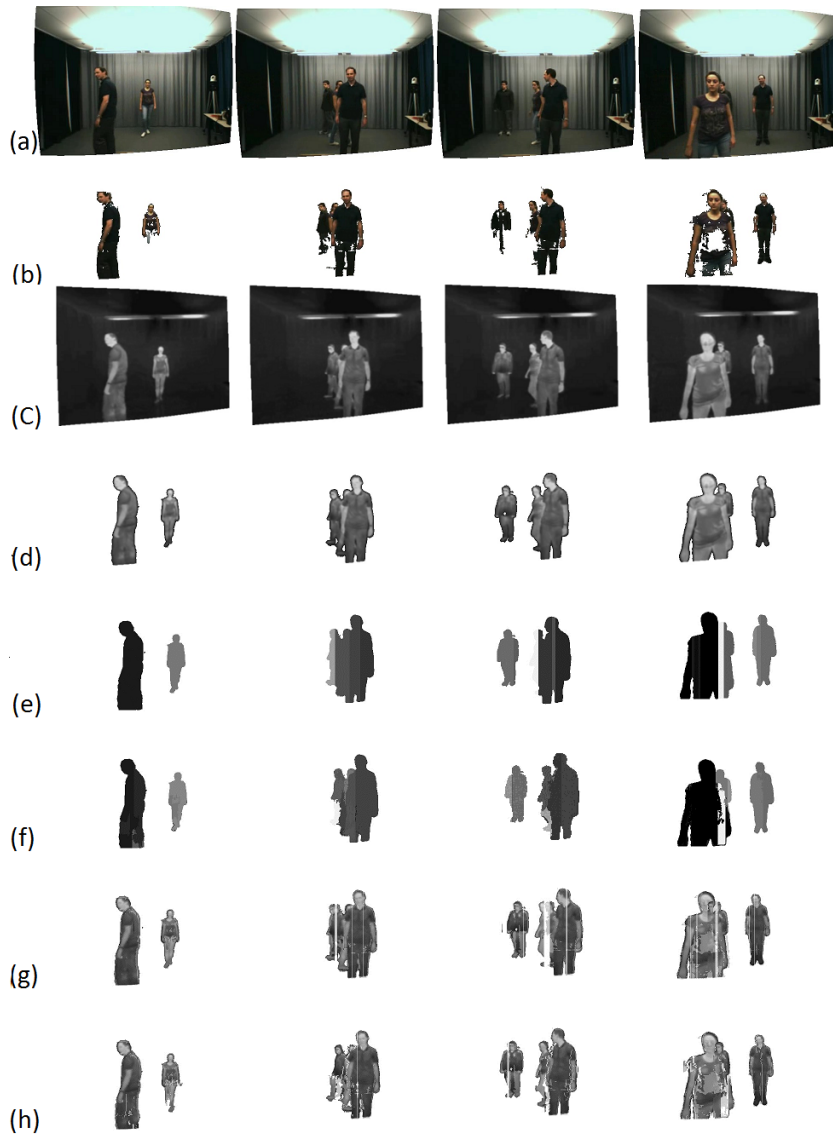


Figure 14: Comparison of MI+DV method in [7] and our proposed method LSS+MS+DV for our summer video using imperfect foreground segmentation (mainly misdetection). (a) visible image, (b) visible foreground segmentation, (c) thermal image, (d) thermal foreground segmentation, (e) MI+DV disparity image, (f) LSS+MS+DV disparity image, (g) MI+DV registration, and (h) LSS+MS+DV registration.

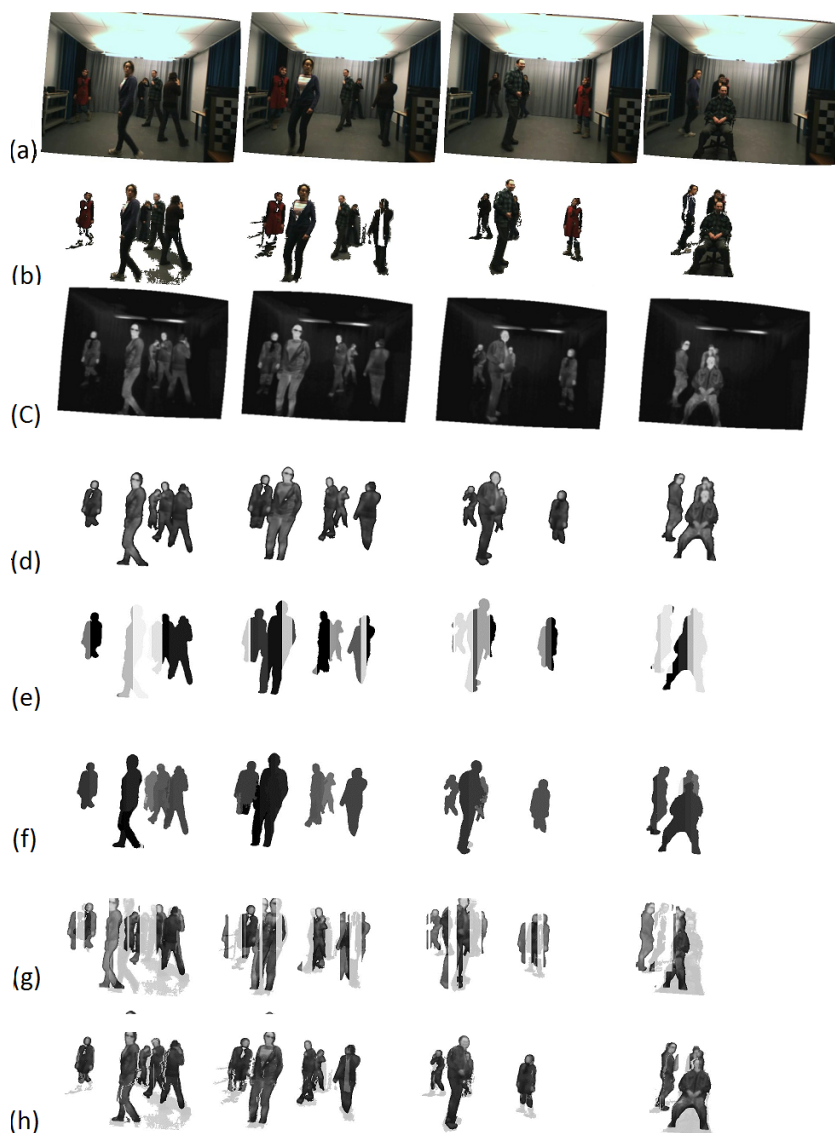


Figure 15: Comparison of MI+DV method in [7] and our proposed method LSS+MS+DV for our winter video using imperfect foreground segmentation (false detection and misdetection). (a) visible image, (b) visible foreground segmentation, (c) thermal image, (d) thermal foreground segmentation, (e) MI+DV disparity image, (f) LSS+MS+DV disparity image, (g) MI+DV registration, and (h) LSS+MS+DV registration.

547 As future direction for this work, we are working on improving the motion
548 segmentation results to obtain more accurate segments and to avoid over
549 segmentation.

550 **References**

- 551 [1] Z. Zhu, T. Huang, Multimodal surveillance: an introduction, in: *Com-*
552 *puter Vision and Pattern Recognition, 2007. CVPR '07. IEEE Confer-*
553 *ence on, 2007*, pp. 1 –6.
- 554 [2] R. Collins, A. Lipton, H. Fujiyoshi, T. Kanade, Algorithms for cooper-
555 ative multisensor surveillance, *Proceedings of the IEEE* 89 (10) (2001)
556 1456 –1477.
- 557 [3] D. Socolinsky, Design and deployment of visible-thermal biometric
558 surveillance systems, in: *Computer Vision and Pattern Recognition,*
559 *2007. CVPR '07. IEEE Conference on, 2007*, pp. 1 –2.
- 560 [4] J. W. Davis, V. Sharma, Background-subtraction using contour-based
561 fusion of thermal and visible imagery, *Comput. Vis. Image Underst.* 106.
- 562 [5] A. Leykin, Thermal-visible video fusion for moving target tracking and
563 pedestrian classification, in: *In Object Tracking and Classification in*
564 *and Beyond the Visible Spectrum Workshop at the International Con-*
565 *ference on Computer Vision and Pattern Recognition, 2007*, pp. 1–8.
- 566 [6] J. Han, B. Bhanu, Fusion of color and infrared video for moving human
567 detection, *Pattern Recognition* 40 (6) (2007) 1771 – 1784.

- 568 [7] S. J. Krotosky, M. M. Trivedi, Mutual information based registration
569 of multimodal stereo videos for person tracking, *Computer Vision and*
570 *Image Understanding* 106 (2-3) (2007) 270 – 287.
- 571 [8] H.-M. Chen, P. Varshney, M.-A. Slamani, On registration of regions
572 of interest (roi) in video sequences, in: *IEEE Conference on Advanced*
573 *Video and Signal Based Surveillance (AVSS 2003)*, 2003, pp. 313 – 318.
- 574 [9] C. Fookes, A. Maeder, S. Sridharan, J. Cook, Multi-spectral stereo image
575 matching using mutual information, in: *3D Data Processing, Visualiza-*
576 *tion and Transmission, 2004. 3DPVT 2004. Proceedings. 2nd Interna-*
577 *tional Symposium on*, 2004, pp. 961 – 968.
- 578 [10] G. Egnal, Mutual information as a stereo correspondence measure, Tech.
579 Rep. MS-CIS-00-20, University of Pennsylvania.
- 580 [11] E. Shechtman, M. Irani, Matching local self-similarities across images
581 and videos, in: *IEEE Conference on Computer Vision and Pattern*
582 *Recognition (CVPR 2007)*, 2007, pp. 1 –8.
- 583 [12] S. Walk, N. Majer, K. Schindler, B. Schiele, New features and insights
584 for pedestrian detection, in: *IEEE Conference on Computer Vision and*
585 *Pattern Recognition (CVPR 2007)*, 2010, pp. 1030–1037.
- 586 [13] J. Yang, Y. Li, Y. Tian, L. Duan, W. Gao, Group-sensitive multiple
587 kernel learning for object categorization, in: *IEEE 12th International*
588 *Conference on Computer Vision (ICCV 2009)*, 2009, pp. 436 –443.
- 589 [14] A. Vedaldi, V. Gulshan, M. Varma, A. Zisserman, Multiple kernels for

- 590 object detection, in: IEEE 12th International Conference on Computer
591 Vision (ICCV 2009), 2009, pp. 606 –613.
- 592 [15] A. Torabi, G.-A. Bilodeau, Local self-similarity as a dense stereo corre-
593 spondence measure for themal-visible video registration, in: Computer
594 Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Com-
595 puter Society Conference on, 2011, pp. 61 –67.
- 596 [16] P. O. Hoyer, P. Dayan, Non-negative matrix factorization with sparse-
597 ness constraints, *Journal of Machine Learning Research* 5 (2004) 1457–
598 1469.
- 599 [17] A. Torabi, M. Najafianrazavi, G. Bilodeau, A comparative evaluation
600 of multimodal dense stereo correspondence measures, in: *Robotic and
601 Sensors Environments (ROSE)*, 2011 IEEE International Symposium on,
602 2011, pp. 143 –148.
- 603 [18] J. Heikkila, O. Silven, A four-step camera calibration procedure with
604 implicit image correction, in: *Computer Vision and Pattern Recognition,
605 1997. Proceedings.*, 1997 IEEE Computer Society Conference on, 1997,
606 pp. 1106 –1112.
- 607 [19] A. Gil, O. M. Mozos, M. Ballesta, O. Reinoso, A comparative evaluation
608 of interest point detectors and local descriptors for visual slam, *Mach.
609 Vision Appl.* 21 (2010) 905–920.
- 610 [20] R. Mayoral, M. Aurnhammer, Evaluation of correspondence errors for
611 stereo, *Pattern Recognition, International Conference on* 4 (2004) 104–
612 107.

- 613 [21] B. Shoushtarian, H. E. Bez, A practical adaptive approach for dynamic
614 background subtraction using an invariant colour model and object
615 tracking, *Pattern Recogn. Lett.* 26 (1) (2005) 5–26.
- 616 [22] A. Ogale, Y. Aloimonos, A roadmap to the integration of early visual
617 modules, *International Journal of Computer Vision* 72 (2007) 9–25.
- 618 [23] D. Comaniciu, P. Meer, Mean shift analysis and applications, in: *The*
619 *Proceedings of the Seventh IEEE International Conference on Computer*
620 *Vision, (ICCV 1999), Vol. 2, 1999, pp. 1197 –1203 vol.2.*