

Thermal-Visible Registration of Human Silhouettes: a Similarity Measure Performance Evaluation

Guillaume-Alexandre Bilodeau^{a,*}, Atousa Torabi^b, Pierre-Luc St-Charles^a,
Dorra Riahi^a

^a*LITIV Lab., Department of Computer and Software Engineering,
École Polytechnique de Montréal,
P.O. Box 6079, Station Centre-ville, Montréal
Québec, Canada, H3C 3A7*
^b*LISA, Dept. IRO,
Université de Montréal,
Montréal, Québec, Canada, H2C 3J7*

Abstract

When dealing with the registration of information from different image sources, the *de facto* similarity measure used is Mutual Information (MI). Although MI gives good performance in many image registration applications, recent works in thermal-visible registration have shown that other similarity measures can give results that are as accurate, if not more than MI. Furthermore, some of these measures also have the advantage of being calculated independently from each image to register, which allows them to be integrated more easily in energy minimization frameworks. In this article, we investigate the accuracy of similarity measures for thermal-visible image registration of human silhouettes, including MI, Sum of Squared Differences

*Corresponding author

Email addresses: gabilodeau@polymtl.ca (Guillaume-Alexandre Bilodeau),
torabi@iro.umontreal.ca (Atousa Torabi), pierre-luc.st-charles@polymtl.ca
(Pierre-Luc St-Charles), dorra.riahi@polymtl.ca (Dorra Riahi)

(SSD), Normalized Cross-Correlation (NCC), Histograms of Oriented Gradients (HOG), Local Self-Similarity (LSS), Scale-Invariant Feature Transform (SIFT), Speeded-Up Robust Features (SURF), Census, Fast Retina Keypoint (FREAK), and Binary Robust Independent Elementary Feature (BRIEF). We tested the various similarity measures in dense stereo matching tasks over 25 000 windows to have statistically significant results. To do so, we created a new dataset in which one to five humans are walking in a scene in various depth planes. Results show that even if MI is a very strong performer, particularly for large regions of interest (ROI), LSS gives better accuracies when ROI are small or segmented into small fragments because of its ability to capture shape. The other tested similarity measures did not give consistently accurate results.

Keywords: Multispectral imagery, similarity measures, dense stereo matching, thermal-visible registration, thermal camera, visible camera

1 Introduction

In the recent years, there has been a growing interest in visual surveillance using multimodal sensors in both civilian and military applications. One of the fundamental issues associated with thermal-visible imagery is the matching and registration of pairs of images captured by the two different types of sensors. Unlike visible sensors that capture reflected light, IR sensors capture thermal radiations reflected and emitted by an object in a scene. Although not very well documented, most similarity measures used for registering visible images are not applicable for thermal-visible image registration because of the differences in imaging characteristics of thermal and visible

11 cameras. Furthermore, scene content at room temperature cannot be regis-
12 tered because it does not convey any textural information. For this reason,
13 most works in thermal-visible imagery focus on hot bodies, like for example,
14 people.

15 For the registration of people, two families of methods exist. First, there
16 are the sparse registration methods. In these methods, people are considered
17 as planar objects that can be registered by just matching some points. Most
18 of these methods are based on matching contour points. Using edges is one
19 of the most popular method as their magnitudes and orientations may match
20 between infrared and visible for some object boundaries [1, 2, 3]. Raw edges
21 alone are not very reliable, so they may be considered as connected groups
22 for correspondence [1, 3]. Other methods use polygonal approximation of
23 people and matches them using their vertices [4, 5]. Although these methods
24 allow fast registration of people, the resulting transformation does not allow
25 to capture fine depth details, like the position of the arms or legs relative to
26 the body.

27 For accounting for depth details, the second family is based on dense cor-
28 respondences between the visible and the thermal human silhouettes. Typ-
29 ically, in that case, registration is performed at every pixel on the human
30 body by comparing the pixels inside a window using a similarity measure.
31 We define as a similarity measure any function that returns a value that
32 indicates a level of similarity. For example, Mutual Information (MI) is one
33 of the most popular similarity measure [6, 7, 8, 9], but recently, Local Self-
34 Similarity (LSS), a local image descriptor, was also proposed for this purpose.
35 A matching window is described as a dense collection of LSS descriptors, and

36 the similarity of two windows is the Euclidean distance between the respec-
37 tive collections of descriptors [10, 11].

38 The proposal of new local image descriptors (LID) is a very active field in
39 computer vision. These LIDs are typically proposed to allow discriminative
40 matches between regions in visible images. For each of these new LIDs, the
41 question is: can they be used to successfully register people in visible and
42 thermal images? This question is worth answering for two reasons: 1) MI
43 does not provide accurate registration all the time [6, 7, 8], and 2) sometime,
44 surprisingly, other LIDs like Local sSelf-Similarity (LSS) have been shown to
45 outperform MI [10, 11]. In this article, our goal is to study the applicability
46 of various LIDs or other similarity measures to the problem of registering
47 people (or any other bodies not at room temperature) in visible and thermal
48 imagery. To test the viability of various similarity measures, we use them in
49 the context of typical windows-based matching, where the potential measures
50 or LIDs are applied over all the pixels in a window to find correspondences.

51 We began some work in that regard in [12]. However, the experiments
52 were more limited (around 300 test windows) and less similarity measures
53 were compared. In this study, we experiment on more similarity measures
54 and we test them on over more than 25 000 windows to have statistically
55 significant results for human silhouette registration. We compare the simi-
56 larity measures both for winner takes all (WTA) sliding window matching
57 and disparity voting (DV).

58 In section 2, we present the similarity measures we tested. Section 3
59 describes the details of our camera setup, dataset, test procedure, and evalu-
60 ation criteria. In section 4, we present and discuss our experimental results.

61 Finally, we conclude the paper with a general discussion in section 5.

62 **2. Tested similarity measures**

63 We tested three broad categories of similarity measures:

- 64 1. Similarity measures that are calculated across pixels of the two win-
65 dows, namely Mutual Information (MI), Sum of Squared Differences
66 (SSD), and Normalized Cross-Correlation (NCC);
- 67 2. Traditional LIDs that model data as distributions, namely Histograms
68 of Oriented Gradients (HOG), Local Self-Similarity (LSS), Scale Invari-
69 ant Feature Transform (SIFT), Speeded-Up Robust Features (SURF).
70 In this case, they are applied as dense collection of features compared
71 with a distance to be used as similarity measure;
- 72 3. LIDs based on binary comparisons of pixels, namely Census, Fast
73 REtinA Keypoint (FREAK), Binary Robust Independent Elementary
74 Feature (BRIEF). In this case they are also applied as dense collections.

75 These similarity measures represent only a subset of possible measures, as
76 any LIDs could be formulated as similarity measures. However, they repre-
77 sent a good sample of measures as they cover the main categories of measures
78 and LIDs and they all show good discriminative power when comparing re-
79 gions in visible images.

80 NCC is a classic similarity measure that consists in a pixel-wise cross-
81 correlation of two image regions normalized by the overall intensity difference
82 [13]. It is defined for two windows on a pair of images as

$$NCC(W_1, W_2) = \frac{\sum_{x,y}(W_1(x, y) - \bar{W}_1) * (W_2(x, y) - \bar{W}_2)}{\sqrt{\sum_{x,y}(W_1(x, y) - \bar{W}_1)^2 * \sum_{x,y}(W_2(x, y) - \bar{W}_2)^2}}, \quad (1)$$

83 where W_1 and W_2 are two matching windows on a pair of thermal and
 84 visible images, and \bar{W}_1 and \bar{W}_2 are the mean pixel intensities in the windows.
 85 This measure relies basically on similar intensity patterns. This is similar for
 86 SSD, which is defined as

$$SSD(W_1, W_2) = \sum_{x,y}(W_1(x, y) - W_2(x, y))^2. \quad (2)$$

87 MI computes the statistical co-occurrence of pixel-wise image patterns
 88 inside a window on pair of images using

$$MI(W_1, W_2) = \sum_{X \in W_1} \sum_{Y \in W_2} P(X, Y) \log \frac{P(X, Y)}{P(X)P(Y)}, \quad (3)$$

89 where $P(X, Y)$, is the joint probability mass function of intensities and
 90 $P(X)$ and $P(Y)$ are the marginal probability functions. $P(X, Y)$ is cal-
 91 culated by creating a two-dimensional histogram that records the number of
 92 co-occurrences of thermal and visible intensity values in W_1 and W_2 . The
 93 probabilities are then obtained by normalizing the histogram by the sum
 94 of the joint histogram entries. The marginal probabilities $P(X)$ and $P(Y)$
 95 are then obtained by summing $P(X, Y)$ over the grayscale or thermal in-
 96 tensities. MI relies on the co-occurrence of patterns that do not need to be
 97 similar. It can match a uniform region with a textured region. This is why
 98 it is successful with multimodal imagery.

99 To use LSS[14], SIFT[15], and SURF[16] as similarity measures, each of
 100 these feature descriptors are calculated densely for all the pixels in a window.

101 For a window of say 10×10 , this gives a collection of 100 descriptors. The
102 similarity measure SM of two windows W_1 and W_2 is given by

$$SM(W_1, W_2) = \sqrt{\sum_{x,y} (f_l(x, y) - f_r(x, y))^2}, \quad (4)$$

103 where f_l and f_r are feature descriptors in two matching windows on a pair of
104 thermal and visible images. SIFT and SURF are gradient-based, while LSS
105 is a local shape descriptor based on the comparison of the self-similarity of a
106 central patch with neighboring patches. As for HOG[17], since it is window-
107 based, we apply it directly on the whole windows. The HOG descriptors are
108 then compared with the Euclidean distance with Eq. 4.

109 Census[18], BRIEF[19] and FREAK[20] are applied the same way as
110 SIFT. They all consist of binary intensity comparisons within a window using
111 different pre-determined patterns. In this case, the binary descriptors of two
112 windows are compared using the Hamming distance to obtain a similarity
113 measure.

114 3. Experimental method

115 We tested the 10 chosen similarity measures on more than 25 000 windows
116 using both a WTA and DV procedure. In this section, we describe our test
117 methodology.

118 3.1. Video Acquisition and Calibration

119 We used synchronized visible-thermal videos of a $5m \times 5m$ room at a
120 fixed temperature of $24^\circ C$ captured by stationary thermal and visible cam-
121 eras with a 12 cm baseline. The scene is relatively close range where different



Figure 1: Calibration images: (a) visible image (b) thermal image.

122 people in different poses and clothing are walking at different depths (between
 123 2-5 meters) from the camera baseline. In order to simplify the stereo match-
 124 ing to a 1D search, we first calibrated the thermal and visible cameras, and
 125 then rectified the images using the intrinsic and extrinsic calibration param-
 126 eters. We used the standard technique available in the camera calibration
 127 toolbox of MATLAB [21]. As for calibration, we placed a checkerboard pat-
 128 tern in front of the cameras. Since in the thermal images, the checkerboard
 129 pattern is not visible at room temperature, we illuminated the scene using
 130 high intensity halogen bulbs placed behind the two cameras. This way, the
 131 dark squares absorb more energy and visually appear brighter than the white
 132 squares. Fig. 1 shows an example of our calibration images.

133 We captured four videos and manually annotated the disparity of people
 134 in 206 frames from extracted foregrounds using [22], which gives a total of
 135 25819 ground-truth point pairs¹. All videos include between 1 to 5 actors
 136 walking around and occluding each other (see Fig. 2). The specifications of
 137 each video are given in table 1. A second subset of ground-truth point pairs

¹Our dataset is publicly available at <http://www.polymtl.ca/litiv/en/vid/index.php>



Figure 2: Vid2cut1, rectified frames 526: (a) visible image (b) thermal image.

Table 1: Test video dataset. GTpp stands for ground-truth point pairs.

	Nb. Frames	Nb. Actors	GTpp for WTA	GTpp for DV
Vid1	4366	1 to 4	11166	1496
Vid2cut1	96	1 to 2	2217	706
Vid2cut2	451	2 to 4	6012	889
Vid3	477	1 to 5	6524	981
Total	5390	1 to 5	25819	3442

138 was selected for experiments with disparity voting that are more computationally expensive.
 139

140 3.2. Test procedures

141 We tested each similarity measure under two matching procedures. The
 142 first consists in finding the best match for a point selected in the visible in the
 143 thermal image by scanning a row in the image column by column (see Fig.
 144 3). This is the winner-takes-all (WTA) procedure and its algorithm is given
 145 in Algo. 1. For this test procedure, we found matches on all the annotated
 146 pixels in the four videos. The annotated pixel coordinate is used as the center

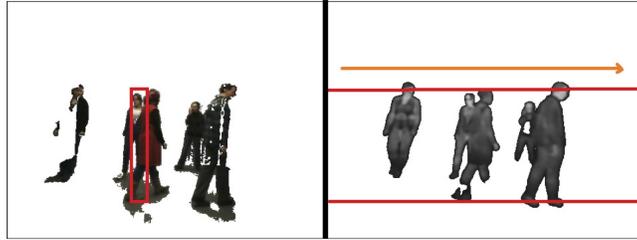


Figure 3: Visible-thermal 1-D sliding window matching.

147 of the matching window. The WTA test procedure was repeated using three
 148 rectangular window sizes of 10×130 (small), 20×130 (medium), and $40 \times$
 149 130 (large) pixels. The heights of the windows are larger than their widths
 150 since our matching target is a human body ROI that has such proportions.
 151 Note that we tested the measures on foreground images that were obtained
 152 using background subtraction [22]. The background subtraction is not perfect
 153 and ROIs might be partially misdetected or some regions might be falsely
 154 detected. For each tested point, we record the column $best_j$ with maximum
 155 similarity value and its associated similarity value $bestSV_j$, in addition to the
 156 similarity value SV_j obtained for all the columns to evaluate the *power of*
 157 *discrimination* (see section 3.3.2). The similarity value is the value resulting
 158 from the application of a similarity measure, for example, the results obtained
 159 from Eq. 1, 2, 3, 4 and the Hamming distance.

160 The second test procedure is based on disparity voting (DV) as used
 161 previously by [6] in the context of thermal-visible registration. In this case,
 162 the disparity at a column is the result of a majority voting of disparities
 163 obtained from neighboring columns (see Algo. 2). The matching process was
 164 repeated using three rectangular window sizes of 10×130 (small), 20×130
 165 (medium), and 40×130 (large) pixels. For our experiments, V in Algo. 2

Algorithm 1 Winner-takes-all testing procedure

```
1: procedure WTA TESTING
2:   NW: number of tested windows
3:   L: image width
4:   for each test window  $W_i$  centered at  $x,y$  in visible image, with  $i=1\dots NW$  do
5:     if testing LID then
6:       Extract LIDs for all pixels in  $W_i$ 
7:     end if
8:     for each window  $W_j$  centered at  $j,y$  in thermal image, with  $j=1\dots L$  do
9:       if testing LID then
10:        Extract LIDs for all pixels in  $W_j$ 
11:        Compute similarity value  $SV_j$  between LIDs in  $W_i$  and  $W_j$ 
12:       else
13:        Compute similarity value  $SV_j$  between  $W_i$  and  $W_j$ 
14:       end if
15:       Save  $best_j$  and  $bestSV_j$  if similarity of  $W_j$  is larger than for previous  $best_j$ 
16:       Save  $SV_j$ 
17:     end for
18:     Compare  $best_j$  with ground-truth
19:     Save disparity error  $e_i$ 
20:   end for
21: end procedure
```

166 was selected as the width of the rectangular window. This procedure was
167 also performed on foreground images. As opposed to the WTA method, the
168 DV method is more robust near occlusion boundaries because neighboring
169 disparities are considered for taking a decision at a particular column.

170 3.3. Evaluation Criteria

171 To evaluate the similarity measures, we have used three metrics. The first
172 two metrics are the *precision* and *recall* while the third metric is the power
173 of discrimination.

174 3.3.1. Precision and Recall

175 We used a criterion based on the number of correct matches for all pairs of
176 tested images as used in [23, 24]. *Precision* and *recall* are defined as follows:

$$177 \textit{precision} = \frac{CM}{MR}, \quad (5)$$

$$\textit{recall} = \frac{CM}{TM}, \quad (6)$$

178 where CM is the number of correct matches, MR is the number of
179 matches retrieved, and TM is the total number of ground-truth matches.
180 In a *precision* versus *recall* curve, a feature with high *recall* and low *preci-*
181 *sion* means that many correct matches as well as many false matches are
182 retrieved. On the other hand, high *precision* and low *recall* means that most
183 matches are correct but many others have been missed. In our experiment,
184 TM is a fixed value that corresponds to the total number of tested windows.

185 Remember that for each tested window, we obtain a $bestSV_j$ value that
186 corresponds to the maximum similarity value obtained when comparing two

Algorithm 2 Disparity voting testing procedure

```
1: procedure DV TESTING
2:   NW: number of tested windows
3:   L: image width
4:   V: Number of windows for voting
5:   for each test window  $W_i$  centered at x,y in visible image, with  $i=1\dots NW$  do
6:     for  $k = -V/2\dots V/2$  do
7:       if testing LID then
8:         Extract LIDs for all pixels in  $W_{i+k}$ 
9:       end if
10:      for each window  $W_j$  centered at j,y in thermal image, with  $j=1\dots L$  do
11:        if testing LID then
12:          Extract LIDs for all pixels in  $W_j$ 
13:          Compute similarity value between LIDs in  $W_{i+k}$  and  $W_j$ 
14:        else
15:          Compute similarity value between  $W_{i+k}$  and  $W_j$ 
16:        end if
17:        Save  $j$  if similarity value of  $W_j$  is larger than for previous  $j$ 
18:      end for
19:      Add a vote for disparity  $d = j - i + k$ 
20:    end for
21:    Select  $d$  with more votes
22:    Compare  $d$  with ground-truth
23:    Save disparity error  $e_i$ 
24:  end for
25: end procedure
```

187 windows with a given approach. Depending on the tested window, the max-
188 imum similarity value can be large or small. Thus, in order to obtain dif-
189 ferent recall-precision points to draw our curves, we sort our results using
190 their $bestSV_j$ values and define MR as the n first matches of the sorted re-
191 sults vector ($n \in [1, TM]$). We then define CM as a subset of MR in which
192 all $best_j$ values have a disparity error smaller than 3 pixels (with respect to
193 the ground truth). Hence, the first points ($n \ll TM$) should present high
194 *precision* (albeit low *recall*) and the latter points ($n \gg 1$) the opposite.

195 3.3.2. power of discrimination

196 To assess the reliability of matches, correct ground-truth overlap is not the
197 only important factor: the matches must also be discriminative. The *power*
198 *of discrimination* metric describes the level of distinctiveness of a match
199 compared to its neighboring matches on the SV_j versus column j curve.
200 In order to evaluate the *power of discrimination* of similarity measures, we
201 study the shape of SV_j over an interval of columns D for all the matches of
202 all pairs of ground-truth points. Note that for some similarity measures, a
203 good match is a minimum, while for others a good match is a maximum. To
204 illustrate the *power of discrimination*, we will consider the former. In that
205 case, a reliable match is located on an isolated minimum on the SV_j versus j
206 curve and has a SV_j value much smaller than its neighboring points. In order
207 to evaluate the isolation of the minimum, the SV_j values computed by the
208 WTA test procedure on the interval of columns D around the minimum are
209 first normalized and sorted increasingly and are transformed to the interval
210 $[0, 1]$ which results in SV'_j . Then, N is computed by counting the number of
211 values SV'_j that are less than a pre-computed small threshold α , ignoring the

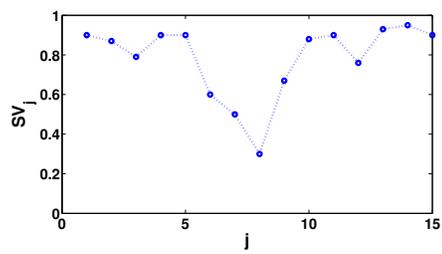
212 minimum. α has the same value for evaluating all similarity measures. Third,
213 a quality measure s (the s value) is computed by dividing N by the size of
214 the interval D . So $s = 0$ corresponds to the most isolated minimum (best
215 performance), and $s = 1$ corresponds to the least isolated minimum. Finally,
216 for each similarity measure, a graph of accumulated frequencies (AF) of the s
217 values of all matches is computed. Therefore, the similarity measure for which
218 AF reaches a higher value at a smaller s value is the most discriminative.
219 Fig. 4 (a) and (b) show an example where the minimum of SV_j is relatively
220 isolated and N is 2, and fig. 4 (c) and (d) show an example where the
221 minimum is not well isolated and N is 8, which results in higher value of s
222 compared to the previous example.

223 4. Experimental Results

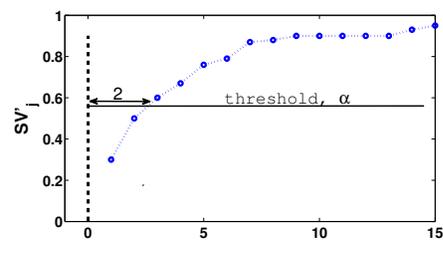
224 We will first present and discuss the results for the WTA test procedure,
225 and we will then present the results for the DV test procedure.

226 4.1. Results with WTA test procedure

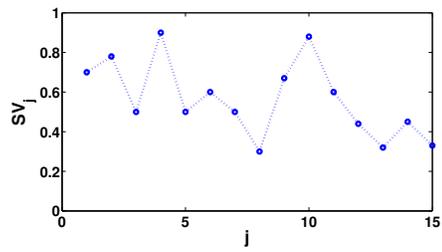
227 Fig. 5, 6, 7 present the *precision-recall* curves obtained for the three
228 tested window sizes. Globally, as expected, a better precision is achieved with
229 a larger window size, as more information is available to measure similarity.
230 The two best performing similarity measures are LSS and MI. LSS tends
231 to be better for smaller windows, while MI tends to be better for larger
232 ones. The key of MI success is that it can consider similar different looking
233 patterns if they co-occur in the matching windows. In the case of LSS, it
234 works well because it is a local shape descriptor that evaluates the similarity
235 of surrounding patches with respect to a center patch. Two shapes can be



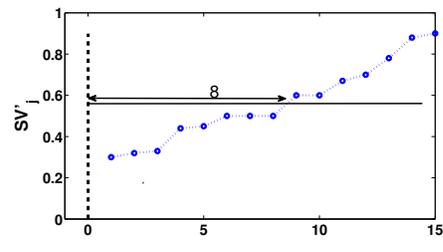
(a)



(b)



(c)



(d)

Figure 4: (a) and (c) Similarity distance SV_j versus column j curve. (b) and (d) Sorted SV'_j curve.

236 considered similar even if they are composed of globally different patterns
237 if laid out the same way. In a sense, the principle behind MI and LSS are
238 similar, which explains their similar performance. On the other hand, our
239 experiments show that pixel by pixel comparisons as performed by SSD and
240 NCC are not successful because the textural content between visible and
241 infrared silhouette is too different. BRIEF and Census suffer from similar
242 problems as they model fine small-scale textural patterns. Interestingly, even
243 though FREAK is also based on binary comparisons, it ranks third or fourth
244 depending on window size. This can be explained by the larger spatial extent
245 of the comparisons done for constructing the FREAK descriptor, and by the
246 fact that these comparisons are done on mean intensity values. As a result,
247 it seems that FREAK is capable of capturing some shape information of the
248 people’s silhouette.

249 Also SURF gives performances that are close to FREAK, but superior to
250 SIFT. These descriptors can sometimes give acceptable results, because they
251 are gradient-based, and gradients were previously shown to be a common
252 feature of visible and thermal human silhouette, particularly at the body’s
253 boundary [1, 2, 3]. The results with HOG are disappointing considering that
254 this feature descriptor is quite similar to SIFT in essence and considering
255 previous results obtained by Torabi *et al.* [12] that were reporting that HOG
256 was performing relatively well. We believe that HOG may not perform as
257 well as SIFT because of intensity normalization that may reduce the saliency
258 of some edges in the thermal images. Furthermore, compared to previously
259 reported results, our new results are coherent with them in the sense that
260 HOG was clearly not as good as LSS and MI, but it seems that using a

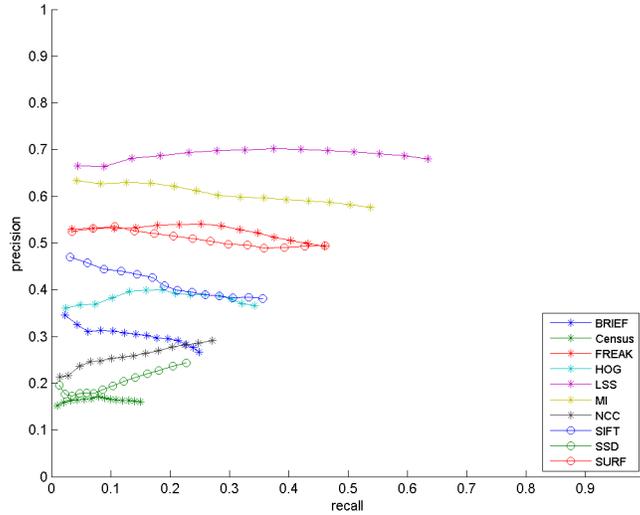


Figure 5: *Precision-recall* curves for small windows (10×130) for the WTA test procedure.

261 much larger number of test points resulted in a significant drop in its relative
 262 performance.

263 We would also like to point out, that the *recall* does not have a large
 264 impact on *precision*. It seems that the similarity measures are distinctive or
 265 not. That is, the absolute similarity value is not indicative of distinctiveness.

266 Let us now consider the *power of discrimination* of our 10 similarity
 267 measures. Fig. 8, 9, 10 present the *power of discrimination* curves obtained
 268 for the three tested window sizes. Recall that these curves allow assessing
 269 the distinctiveness of the best match compared to the following best matches;
 270 that is, it describes how isolated the best match is, with a match being more
 271 isolated when $s = 0$. Thus, a more discriminative measure should have
 272 more correct matches with small s value. These curves show that MI has

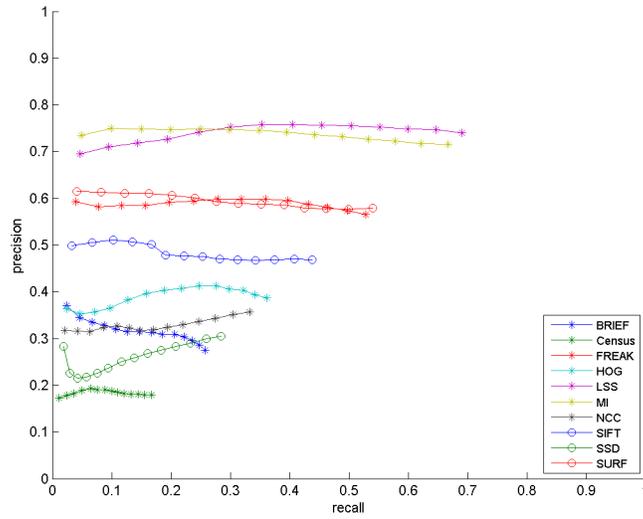


Figure 6: *Precision-recall* curves for medium windows (20×130) for the WTA test procedure.

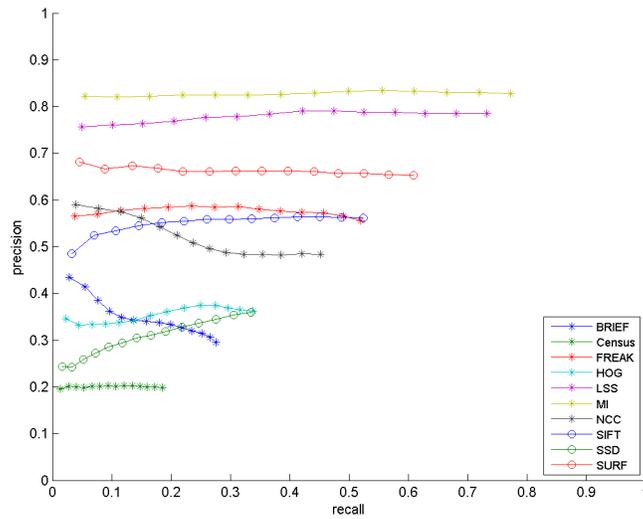


Figure 7: *Precision-recall* curves for large windows (40×130) for the WTA test procedure.

273 the most discriminative matches, followed by LSS and NCC. There is not
274 much to add to the discussion about LSS and MI, except that MI seems to
275 have a small advantage for distinctiveness, particularly with small windows.
276 This contradicts to some extent the results of the *precision-recall* curves that
277 show that LSS is better than MI for small windows. It should be interpreted
278 as LSS gives plenty of matches for smaller window, but the good matches
279 obtained by MI are from more isolated maximums (i.e. more certain). The
280 results from NCC can be viewed a bit as a surprise, but what it tells is that
281 NCC does not work in general for such an application, but when it gives a
282 good match, this good match is very discriminative. Except for Census that
283 is not at all discriminative, all other measures are more or less similar for
284 their *power of discrimination*.

285 4.2. Results with DV test procedure

286 For the DV test procedure, we just report the *recall* because our exper-
287 iments just give us a disparity for each point. Fig. 11, 12, 13 present the
288 *recall* obtained for the three tested window sizes. Globally, using DV does
289 not benefit the best similarity measures. The *recall* for MI, FREAK, and
290 LSS does not improve much over the WTA procedure. This is because they
291 are already quite discriminative. However, the voting mechanism helps most
292 of the other methods. This is particularly true for BRIEF and SSD. BRIEF
293 improves at the level of SURF and SIFT. This indicates that even if disparity
294 voting will not systematically improve registration accuracy, it can improve
295 robustness for less effective similarity measures. Concerning MI and LSS, we
296 still get the same conclusions. LSS is a better for small windows, while MI
297 is better for larger windows. Both methods are about on par for medium

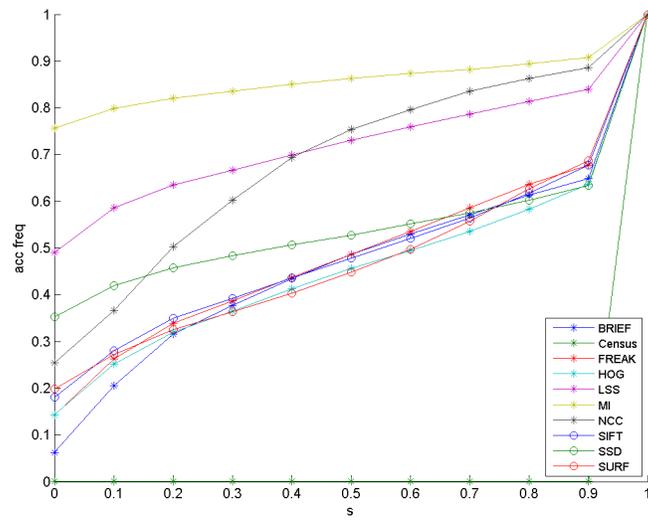


Figure 8: *Power of discrimination* curves for small windows (10×130) for the WTA test procedure.

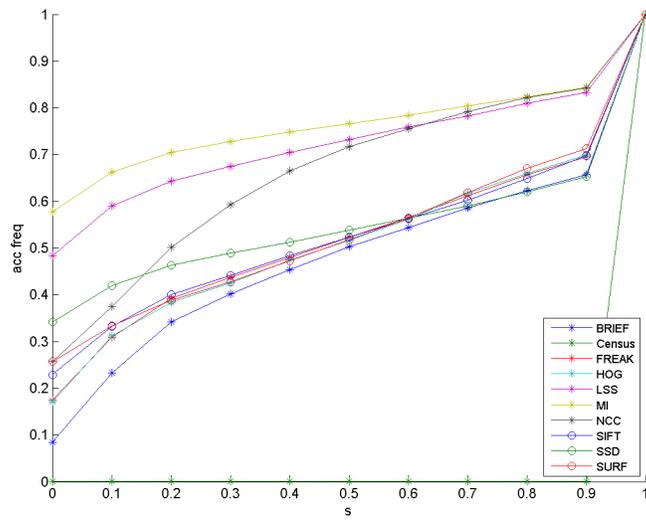


Figure 9: *Power of discrimination* curves for medium windows (20×130) for the WTA test procedure.

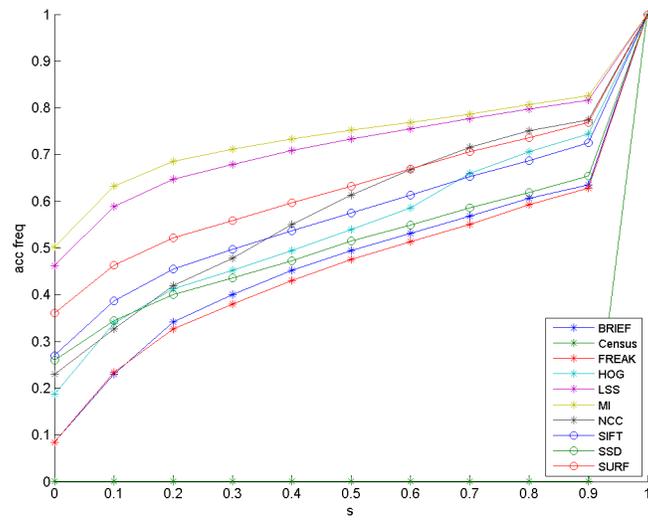


Figure 10: *Power of discrimination* curves for large windows (40×130) for the WTA test procedure.

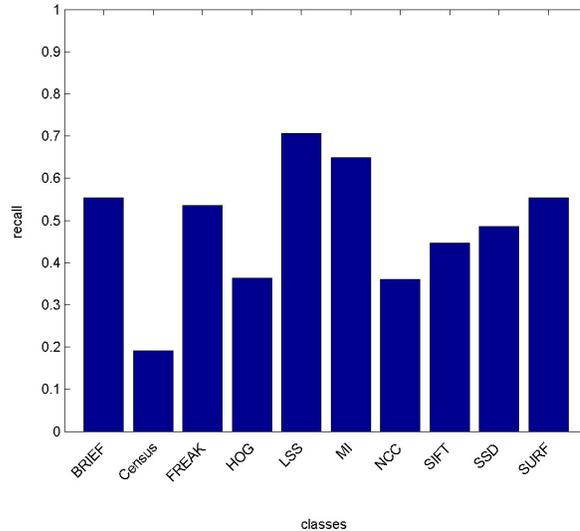


Figure 11: *Recall* for small windows (10×130) for the WTA test procedure.

298 windows. Again, globally, for larger window sizes, the results are better be-
 299 cause more information is available. By more information, we mean that
 300 a larger window will necessarily include more boundary information that is
 301 very discriminative.

302 5. General discussion and conclusions

303 From our results, we observed that although MI is certainly a good choice
 304 for registering human silhouettes, it is not always the most accurate similarity
 305 measure for thermal-visible registration. As such, we should be attentive to
 306 new feature descriptors that may help us obtain future results that are even
 307 more accurate than what we have now. From our tests, we can conclude
 308 that:

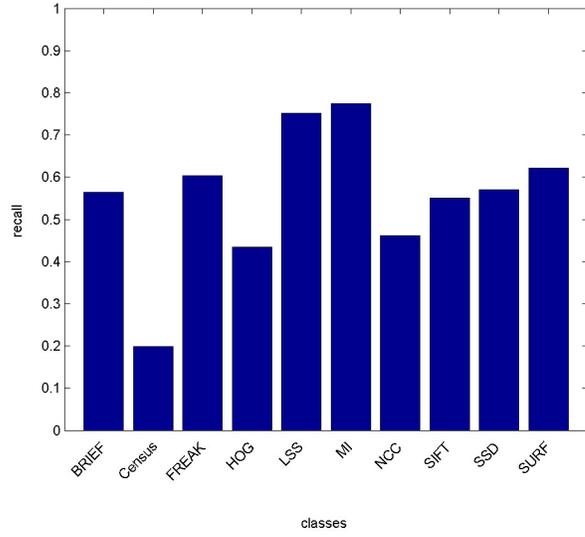


Figure 12: *Recall* for medium windows (20×130) for the WTA test procedure.

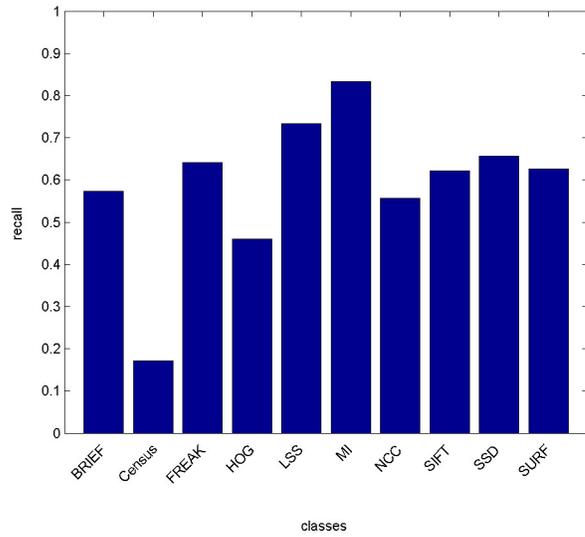


Figure 13: *Recall* for large windows (40×130) for the WTA test procedure.

- 309 • to obtain an accurate registration of thermal and visible silhouettes, the
310 similarity measure chosen should be designed to account principally for
311 the layout of patterns. This is the case for MI and LSS;
- 312 • similarity measures that model fine small-scale texture cannot work in
313 this context;
- 314 • if the objects of interest are small or segmented into small fragments,
315 or if there are many occlusions between objects which require small
316 windows for treatment, LSS is the best choice to obtain accurate results;
- 317 • if the objects are large, MI is still the best similarity measure;
- 318 • *precision* and *recall* metrics are more useful metrics than the *power of*
319 *discrimination* to evaluate a new similarity measure. This was illus-
320 trated by the results for NCC;
- 321 • disparity voting improves the robustness of the similarity measures, but
322 the effect is more obvious for weaker similarity measures.

323 6. Acknowledgements

324 This research was supported by a Natural Sciences and Engineering Re-
325 search Council of Canada (NSERC) Discovery grant No. 311869-2010.

- 326 [1] E. Coiras, J. Santamaria, C. Miravet, Segment-based registration tech-
327 nique for visual-infrared images, *Optical Engineering* 39 (2000) 282–289.
- 328 [2] M. I. Elbakary, M. K. Sundareshan, Multi-modal image registration
329 using local frequency representation and computer-aided design (cad)
330 models, *Image Vision Computing* 25 (5) (2007) 663–670.

- 331 [3] Z. Zhu, T. Huang, Multimodal surveillance: an introduction, in: Com-
332 puter Vision and Pattern Recognition, 2007. CVPR '07. IEEE Confer-
333 ence on, 2007, pp. 1 –6.
- 334 [4] G. A. Bilodeau, P. St-Onge, R. Garnier, Silhouette-based features for
335 visible-infrared registration, in: Computer Vision and Pattern Recog-
336 nition Workshops (CVPRW): Object Tracking and Classification Be-
337 yond the Visible Spectrum, 2011 IEEE Computer Society Conference
338 on, 2011, pp. 68–73.
- 339 [5] S. Sonn, G.-A. Bilodeau, P. Galinier, Fast and accurate registration of
340 visible and infrared videos, in: Computer Vision and Pattern Recog-
341 nition Workshops (CVPRW): Perception Beyond the Visible Spectrum,
342 2013 IEEE Conference on, 2013, pp. 308–313.
- 343 [6] S. J. Krotosky, M. M. Trivedi, Mutual information based registration
344 of multimodal stereo videos for person tracking, *Computer Vision and*
345 *Image Understanding* 106 (2-3) (2007) 270 – 287.
- 346 [7] H.-M. Chen, P. Varshney, M.-A. Slamani, On registration of regions
347 of interest (roi) in video sequences, in: *IEEE Conference on Advanced*
348 *Video and Signal Based Surveillance (AVSS 2003)*, 2003, pp. 313 – 318.
- 349 [8] C. Fookes, A. Maeder, S. Sridharan, J. Cook, Multi-spectral stereo image
350 matching using mutual information, in: *3D Data Processing, Visualiza-*
351 *tion and Transmission (3DPVT 2004)*. Proceedings. 2nd International
352 Symposium on, 2004, pp. 961 – 968.

- 353 [9] S. K. Kyoung, H. L. Jae, B. R. Jong, Robust multi-sensor image reg-
354 istration by enhancing statistical correlation, in: Information Fusion,
355 2005 8th International Conference on, Vol. 1, 2005, p. 7.
- 356 [10] A. Torabi, G.-A. Bilodeau, Local self-similarity-based registration of
357 human {ROIs} in pairs of stereo thermal-visible videos, Pattern Recog-
358 nition 46 (2) (2013) 578 – 589.
- 359 [11] A. Torabi, G.-A. Bilodeau, A lss-based registration of stereo thermalvis-
360 ible videos of multiple people using belief propagation, Computer Vision
361 and Image Understanding 117 (12) (2013) 1736 – 1747.
- 362 [12] A. Torabi, M. Najafianrazavi, G. Bilodeau, A comparative evaluation
363 of multimodal dense stereo correspondence measures, in: Robotic and
364 Sensors Environments (ROSE), 2011 IEEE International Symposium on,
365 2011, pp. 143 –148.
- 366 [13] J. Sarvaiya, S. Patnaik, S. Bombaywala, Image registration by template
367 matching using normalized cross-correlation, in: Advances in Comput-
368 ing, Control, Telecommunication Technologies (ACT '09). International
369 Conference on, 2009, pp. 819 –822.
- 370 [14] E. Shechtman, M. Irani, Matching local self-similarities across images
371 and videos, in: IEEE Conference on Computer Vision and Pattern
372 Recognition (CVPR 2007), 2007, pp. 1 –8.
- 373 [15] D. G. Lowe, Distinctive image features from scale-invariant keypoints,
374 International Journal of Computer Vision 60 (2) (2004) 91–110.

- 375 [16] H. Bay, T. Tuytelaars, L. Gool, Surf: Speeded up robust features, in:
376 A. Leonardis, H. Bischof, A. Pinz (Eds.), Computer Vision ECCV 2006,
377 Vol. 3951 of Lecture Notes in Computer Science, Springer Berlin Hei-
378 delberg, 2006, pp. 404–417.
- 379 [17] N. Dalal, B. Triggs, Histograms of oriented gradients for human detec-
380 tion, in: Proceedings of the 2005 IEEE Computer Society Conference
381 on Computer Vision and Pattern Recognition (CVPR'05) - Volume 1 -
382 Volume 01, CVPR '05, 2005, pp. 886–893.
- 383 [18] R. Zabih, J. Woodfill, Non-parametric local transforms for computing vi-
384 sual correspondence, in: J.-O. Eklundh (Ed.), Computer Vision ECCV
385 '94, Vol. 801 of Lecture Notes in Computer Science, Springer Berlin
386 Heidelberg, 1994, pp. 151–158.
- 387 [19] M. Calonder, V. Lepetit, C. Strecha, P. Fua, Brief: Binary robust inde-
388 pendent elementary features, in: K. Daniilidis, P. Maragos, N. Paragios
389 (Eds.), Computer Vision ECCV 2010, Vol. 6314 of Lecture Notes in
390 Computer Science, Springer Berlin Heidelberg, 2010, pp. 778–792.
- 391 [20] A. Alahi, R. Ortiz, P. Vandergheynst, Freak: Fast retina keypoint, in:
392 Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Confer-
393 ence on, 2012, pp. 510–517.
- 394 [21] J.-Y. Bouguet, Camera calibration toolbox for matlab,
395 http://www.vision.caltech.edu/bouguetj/calib_doc/.
- 396 [22] B. Shoushtarian, H. E. Bez, A practical adaptive approach for dynamic

- 397 background subtraction using an invariant colour model and object
398 tracking, *Pattern Recognition Letters* 26 (1) (2005) 5–26.
- 399 [23] A. Gil, O. M. Mozos, M. Ballesta, O. Reinoso, A comparative evaluation
400 of interest point detectors and local descriptors for visual slam, *Machine*
401 *Vision and Applications* 21 (2010) 905–920.
- 402 [24] K. Mikolajczyk, C. Schmid, A performance evaluation of local descrip-
403 tors, *Pattern Analysis and Machine Intelligence*, *IEEE Transactions on*
404 *27* (10) (2005) 1615 –1630.