

Mosaicing of Camera-captured Document Images

1 Jian Liang^a, Daniel DeMenthon^b, David Doermann^b

2 ^a*Jian Liang is with Amazon.com; Seattle, WA; USA.*

3 ^b*Daniel DeMenthon and David Doermann are with University of Maryland;*
4 *College Park, MD; USA.*

5 **Abstract**

6 In this paper we present a method for composing document mosaics from camera-
7 captured images. We decompose the complexity of solving the 8-dof transformation
8 between image pairs into two problems, that is, rectification and registration. This
9 is achievable under a key assumption that sufficient text content forms orthogonal
10 texture flows on the document surface. First, perspective distortion and rotation are
11 removed from images using the texture flow information. Next, the translation and
12 scaling are resolved by a Hough transformation-like voting method. In the image
13 composition part, our contribution is a sharpness based selection process which
14 composes a seamless and blur free mosaic for text content. Experiments show that
15 our approach can produce an accurate, sharp, and high resolution mosaic of a full
16 document page from small image patches captured by a camera with various zooms
17 and poses.

18 *Key words:* Camera-based document analysis, image mosaicing, image registration

19 **1 Introduction**

20 Digital image mosaicing has been studied for several decades, starting from
21 the mosaicing of aerial and satellite pictures, and now expanding into the
22 consumer market for panoramic picture generation. Its success depends on
23 two key components: image registration and image blending. The first aims at
24 finding the geometric relationship between the to-be-mosaiced images, while
25 the latter is concerned with creating a seamless composition.

Email address: jliang@amazon.com (Jian Liang).

26 Many researchers have developed techniques for the special case of document
27 image mosaicing [3,7,8,12,14,15,11,10]. The basic idea is to create a full, *frontal*
28 view of a document page, often too large to capture during a single scan or in
29 a single frame, by stitching together many small patches.

30 If the small images are obtained through flatbed scanners [3,12], image reg-
31 istration is less challenging because the overlapping part of two images differ
32 only by a 2D Euclidean translation (plus slight rotation, if any), and there is
33 no perspective distortion. When cameras are used, it is still possible to work
34 within similar conditions. In some of the reported work the user is simply asked
35 to point the camera straight at the document plane [7,14]. Others reinforce
36 it with hardware support. For example, Nakao et al. [8] attach a downward
37 looking video camera to a mouse such that displacements among images can
38 be derived from mouse movement. In [15] Zappala et al. fix a downward look-
39 ing camera overhead, and move a document on the desktop, which essentially
40 mimics a scanner. The main weakness of either scanners or fixed cameras is
41 their poor portability.

42 With portable cameras, perspective distortion may exist in the images. Regis-
43 tration is still possible. For example, feature point matching is a common ap-
44 proach in general image registration that is robust against projective transfor-
45 mation. There are affine invariant feature point detectors specially designed for
46 text documents [13]. However, registration by itself does not remove the pro-
47 jectivity. Usually, for document mosaicing, the perspective distortion should
48 be removed.

49 Motivated by structure-from-motion methods, Sato et al. moved from still
50 image cameras to video cameras [11,10]. Their prototype system has an on-

51 line stage which tracks feature points across frames and generates a mosaic
52 preview, and an off-line stage which refines the 3D reconstruction and final
53 mosaic. The on-line stage essentially estimates the extrinsic camera parame-
54 ters, i.e., pose or projectivity; the intrinsic parameters are irrelevant, as long
55 as they are constant. In practice, this translates into using a fixed zoom. One
56 disadvantage of video cameras in terms of document mosaicing is their limited
57 resolution and motion blur.

58 Figure 1 shows two patches of a document captured by a camera with differ-
59 ent zoom and poses. These two images differ in their perspective, resolution,
60 brightness, contrast, and sharpness. Although many methods have been pro-
61 posed for image registration ([9,4], to name a few), samples in Figure 1 still are
62 challenging because of large displacement, small overlap, significant perspec-
63 tive distortion, and periodicity of printed text which presents indistinguishable
64 texture patterns everywhere. For example, we were unable to get any mean-
65 ingful results from global registration technique such as the Fourier-Mellin
66 transform [9] whenever the overlapping area accounts for less than one fourth
67 of each image. In terms of local feature point detection, we tested a general
68 detector (PCA-SIFT [4]) with two robust estimators (Graduated Assignment
69 [2] and RANSAC). The result is unsatisfactory because of the large number
70 of outliers (above 90%) in the result from PCA-SIFT. Figure 2 shows the
71 matched feature points found by PCA-SIFT for three pairs of images.

72 For example, in Figure 2(a), where PCA-SIFT is applied to two outdoor
73 scenery images, most of the matches are correct, as shown by the fact that
74 their connection lines have roughly the same length and direction. There are
75 only a few incorrect matches and they stand out clearly. Figure 2(b) shows
76 two document image patches with the same displacement and scale difference.

77 However, the percentage of incorrect matches is significantly higher because
78 the periodicity of text lines and characters makes feature points less distin-
79 guishable from one another. Figure 2(c) shows the matched points between
80 the two images in Figure 1. In this case, the incorrect matches are so over-
81 whelming that it is very difficult to identify any good matches at a glance.
82 Overall, this example shows that while it is easy to locate feature points in
83 document images, it is more difficult to find good matches under perspective
84 distortion and with small overlapping areas.

85 Our goal is to handle images such as those in Figure 1. Our method removes
86 perspective distortion and registers images. While it is possible to first reg-
87 ister images, then remove perspectivity, we found that once perspectivity is
88 first removed, registration becomes easier. In order to estimate 3D structure
89 information and then to remove perspective, a key assumption is that the
90 document consists of sufficient text content which forms two orthogonal tex-
91 ture flows on the surface. In a certain sense, it is a structure-from-texture (or
92 texture flow) method. First we remove perspective distortion and rotation in
93 each image using the orthogonal texture flows. This step leaves a 3-dof trans-
94 formation (a translation and a scaling) between any two overlapping views.
95 Next we find feature point matches using PCA-SIFT. Although outliers still
96 dominate, we are able to filter them out efficiently with a Hough transform-
97 like voting method. After cross-correlation block matching, we obtain a refined
98 registration between two images, where the perspective distortion is already
99 removed.

100 With respect to image blending, there are three possible problems that have
101 not been well addressed for document images. Conventional blending computes
102 the weighted average in an overlapped area, i.e., $f = a_1f_1 + a_2f_2$, where f_1 and

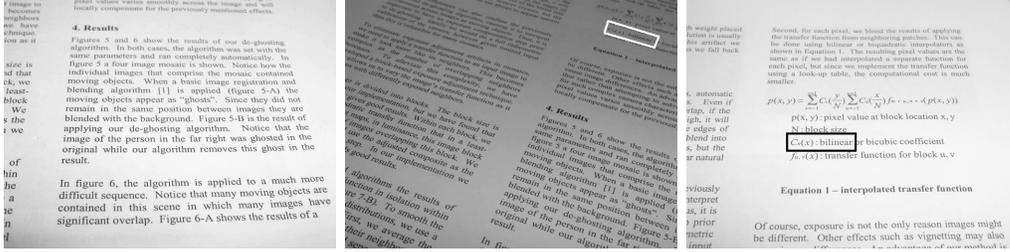


Fig. 1. Image patches of the same document captured from various positions. The same word in two views is marked by overlaid boxes.

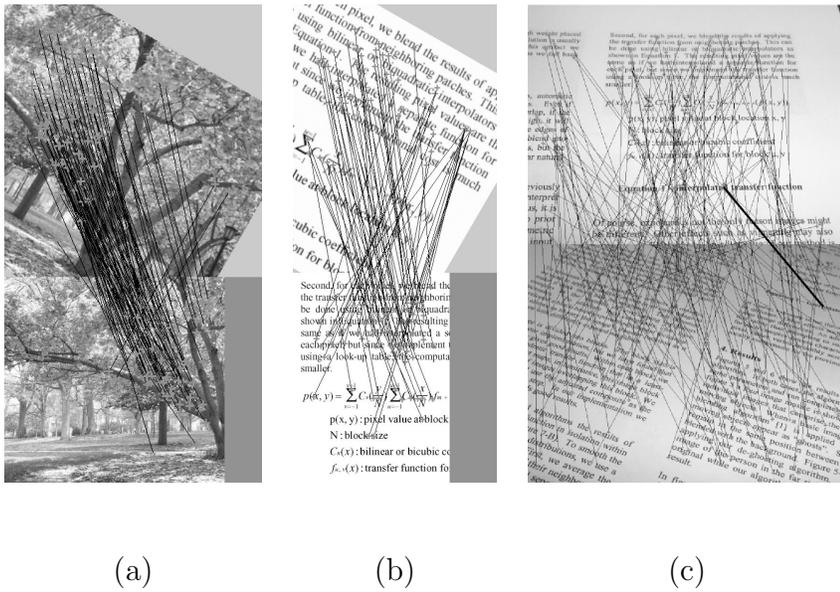


Fig. 2. Match points found by PCA-SIFT. (a) Two sub-images with different scale and rotation generated from a scenery image. (b) Two sub-images from a document image with the same scaling and rotation as in (a). (c) Two camera-captured images of a document page. A thick black line shows one correct match.

103 f_2 are pixel values from two images, a_1 and a_2 are two weights that sum to 1.
 104 By varying the weights, one achieves a gradual transition from one image to
 105 another across the overlapping area. Other more sophisticated methods exist,
 106 which are essentially variations of weighted averaging [1]. Though averaging
 107 usually works well for general images, it is not optimal for document images.
 108 First, the averaging methods treat only the overlapping area. They do not

109 address the overall uneven lighting across images. Second, registration may
 110 have errors. In mis-registered areas, weighted averaging would result in so-
 111 called ‘ghost’ images. Third, two images may have different sharpness because
 112 of different resolution, noise level, zooming, out-of-focus blur, motion blur,
 113 or lighting change. Weighted averaging essentially reduces the sharpness of
 114 the sharper image by blending a blurred image into it. Figure 3 shows the
 115 shortcomings of averaging method. For general scenery or portrait images, a
 116 certain amount of lighting variation and blurring is acceptable and ‘ghosts’
 117 can be softened by blurring. However, for document images, viewers and OCR
 118 algorithms expect sharp contrast between text and background and minimum
 119 lighting variation. Therefore, averaging does not represent the optimal way of
 120 creating document mosaics.

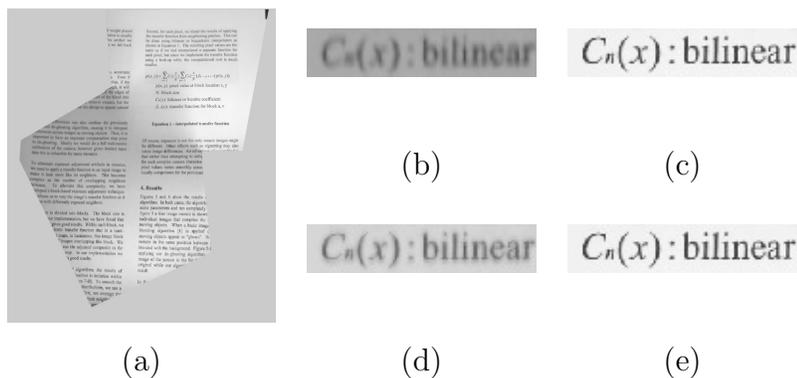


Fig. 3. Challenges for seamless image blending. (a) For two document patches with uneven lighting, their weighted average results in inconsistent contrast across the composite image. (b) A small portion of the darker image. (c) The same portion of the lighter image. (d) Weighted averaging result of (b) and (c) extracted from (a). (e) Our selective image blending result.

121 We treat the inconsistency of lighting by localized histogram normalization,
 122 which balances the brightness and contrast across two images as well as within
 123 each image. Then in the overlapped area, we perform a component level selec-
 124 tive image composition which preserves the sharpness of the printed markings,

125 and ensures a smooth transition near the overlapping area border.

126 A shorter version of our work has appeared in [5]. In this paper we present our
127 method in full details and provide more experimental results. Our prototype
128 system can be illustrated by the pseudo-code in Figure 4. In the next sections
129 we describe the three steps in details.

```
1  Input: two camera-captured document images,  $A_0$  and  $B_0$ 
2  Output: mosaic  $J_1$  free of perspective and rotation
STEP 1: GEOMETRIC RECTIFICATION
3  Detect directions of text lines and vertical strokes;
4  Compute the vanishing points of text lines and vertical strokes;
5  Compute the homography from the vanishing points;
6  Remove the perspective and rotation using the homography:  $A_0 \rightarrow A_1, B_0 \rightarrow B_1$ , where  $A_1$  and  $B_1$  are
   free of perspective and rotation;
STEP 2: IMAGE REGISTRATION
7  Adjust the contrast by local histogram normalization:  $A_1 \rightarrow A_2, B_1 \rightarrow B_2$ ;
8  Find feature point matches using PCA-SIFT:  $(A_2, B_2) \rightarrow M_0$ , where  $M_0$  is a set of matched points;
9  Find the correct matches  $M_1 (\in M_0)$  and the scale  $r$  between  $A_2$  and  $B_2$  using compatible group
   voting:  $M_0 \rightarrow (M_1, r)$ ;
10 Scale  $B_2$  by  $r$ :  $B_2 \rightarrow B'_2$ ;
11 Based on  $M_1$ , compute an initial registration  $H_0$  between  $A_2$  and  $B'_2$ ;
12 Using  $H_0$ , find a set of dense matched points  $M_2$  between  $A_2$  and  $B'_2$  using cross-correlation matching.
13 Using  $M_2$ , compute the final registration  $H_1$  between  $A_2$  and  $B'_2$ ;
STEP 3: SEAMLESS COMPOSITION
14 Compute the “sharpness” maps of  $A_2$  and  $B'_2$ :  $A_2 \rightarrow S_A, B'_2 \rightarrow S_B$ , where  $S_A$  holds the sharpness
   measure of each pixel in  $A_2$ , and so is  $S_B$ ;
15 Using  $H_1$ , composite an initial mosaic  $J_0$  of  $A_2$  and  $B'_2$  by averaging their overlapping part;
16 Find connected component set  $C$  in the binarized version of  $J_0$ ;
17 For each element  $c$  of  $C$  do:
18   Sum up the “sharpness” of all pixels in  $c$ :  $S_A \rightarrow s_A^c, S_B \rightarrow s_B^c$ ;
19   If  $s_A^c > s_B^c$ , copy all pixels in  $c$  from  $A_2$  to final mosaic  $J_1$ ; otherwise, copy from  $B'_2$ ;
20 End for
21 Fill other parts of  $J_1$  by averaging  $A_2$  and  $B'_2$ ;
```

Fig. 4. Workflow of procedure for mosaicing camera-captured documents

130 2 Document Image Rectification

131 Our mosaicing approach removes perspective distortion and rotation in docu-
132 ment images in a step called geometric rectification. This step was described
133 in great detail in [6]. Here we only provide a brief description for completeness.
134 Please refer to [6] for implementation details. First, we detect text line direc-
135 tion and vertical character stroke direction in the image, using local projection
136 profile analysis and directional filter respectively. Figure 2 shows the detected
137 text line and vertical stroke directions superposed on the original text. Sec-
138 ond, we find the vanishing points of these two groups of orthogonal directions
139 using SVD decomposition. From these vanishing points we can estimate the
140 focal length of the camera, the document plane orientation, and finally the
141 homography that maps the two vanishing points to infinity at east and north.
142 The result of geometric rectification is a document image that is free of per-
143 spective distortion and is rotated so that all text lines are horizontal. See [6]
144 for details.

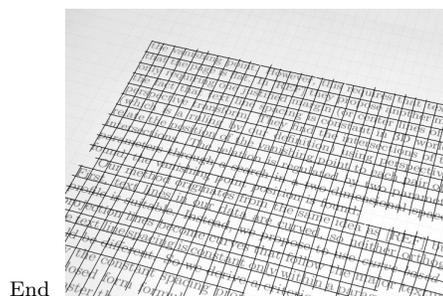


Fig. 5. Text line and vertical stroke directions found in an document image

145 Figure 6 shows how the two rightmost images in Figure 1 are transformed
146 by the rectification (followed by a local histogram normalization described in
147 Section 4).

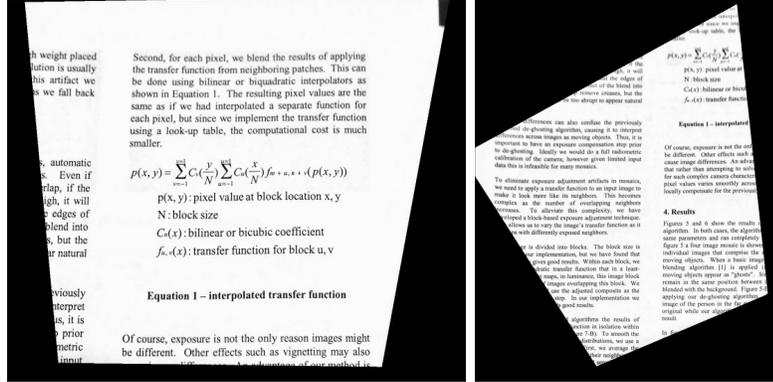


Fig. 6. Document patches after rectification and local histogram normalization.

148 3 Document Image Registration

149 Although projectivity has been removed after geometric rectification, small
 150 overlap, large displacement, and periodicity of texture are still challenging for
 151 common registration methods. For example, the Fourier-Mellin registration
 152 still fails because of insufficient overlapping, and PCA-SIFT still gives a lot of
 153 false matches that defeat Graduated Assignment and make RANSAC ineffec-
 154 tive. However, we are able to filter out the outliers using a Hough transform-
 155 like voting mechanism, since we know only a translation and a scaling remain
 156 to be found.

157 First, let us assume the scale is known. Suppose two images (called A and B)
 158 are placed within the same coordinate system after proper scaling, and the
 159 true translation of image B with respect to image A is (x_0, y_0) . Let $\{p_i\}_{i=1}^N$ be
 160 the feature points in image A, and $\{q_i\}_{i=1}^N$ be the matched points in image B.
 161 If p_i and q_i are a correct match, we have $q_i - p_i = (x_0, y_0)$, and an inequality
 162 otherwise. We compute all the displacements between matched points, i.e., let
 163 $q_i - p_i = (x_i, y_i)$. We have $(x_j, y_j) = (x_k, y_k)$ (we say that they are *compatible*),
 164 where j and k denote any two correct matches. Meanwhile, the probability

165 of having $(x_s, y_s) = (x_t, y_t)$, where either s or t denotes an incorrect match,
166 is extremely low assuming incorrect matches are randomly distributed across
167 the image. We group the matches with equal displacement (within a certain
168 quantization bound) into compatible groups. Ideally, all correct matches are
169 assigned to one group, while each incorrect match constitutes a group of its
170 own. Hence, the correct matches are the matches in the largest group, and
171 their displacements represent the correct translation. In practice, due to the
172 quantization in the histogram used in our compatibility test (see below), some
173 incorrect matches that have similar displacements may be placed in the same
174 group. Even so, the sizes of such groups are highly unlikely to exceed the size
175 of the group of correct matches.

176 If the scale estimate deviates from the correct value, the compatibility mea-
177 sure among correct matches will degrade. A small scale error can be absorbed
178 by the histogram quantization. As the error increases, the group of correct
179 matches will eventually split. Given a completely incorrect scale, the displace-
180 ment distribution of correct matches will be as random as incorrect matches, so
181 the largest compatible group will split into single-match groups. In summary,
182 the largest compatible group is generated when the scale is correct.

183 Based on the above analysis, searching for the largest compatible group of
184 matches as a function of scale can simultaneously solve the problems of finding
185 1) the correct matches, 2) the correct scale, and 3) the correct translation
186 between two images. The specific procedure is as follows:

- 187 (1) For every scale s in a quantized range, construct the compatible groups
188 and let $g(s)$ be the largest.
- 189 (2) Select s^* which maximize $|g(s)|$ and s^* is the correct scale.

190 (3) Find all matches in $g(s^*)$, and compute the mean of their displacements,
191 which is the correct translation.

192 The scale range is quantized on a logarithmic scale. For a given scale, we use
193 a 2D histogram of the match displacements in x and y to find the compatible
194 groups. We divide the 2D displacement space into bins, and the displacement
195 of each match falls in one bin. To address quantization error at bin boundaries,
196 we smooth the 2D histogram by a 3×3 averaging kernel. Then, the bin with
197 the most votes is the largest compatible group. The optimal bin size should
198 be proportional to the average position error of the correctly matched feature
199 points. We use an empirical value, i.e., $1/20$ of the image diagonal length as
200 the bin size. In practice, we find that the sensitivity of the method to this
201 parameter is low (see below).

202 We use PCA-SIFT to find the matches between the two images in Figure 6.
203 Figure 7 shows the number of matches for the first and second largest compat-
204 ible groups found in 2D histograms as a function of the scale. The highest peak
205 in the solid curve identifies the correct scale. At the correct scale, the second
206 largest group (only three votes) is much smaller than the largest group (12
207 votes). This shows good aggregation of correct matches. After examination,
208 we found the second largest group resides in a neighboring bin of the largest
209 group, and the three matches are approximately correct. These two groups
210 would merge if the bin size is increased. With different bin sizes we obtain
211 curves slightly different from those in Figure 7. The correct scale is always
212 found.

213 The figure also shows that when the scale is set slightly larger than the best
214 value, the solid curve drops while the dotted line climbs. This means some

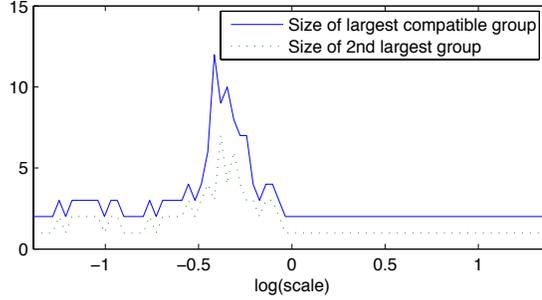


Fig. 7. 2D histogram peak values vs. scales

215 matches in the largest group shift to the second largest group in the neigh-
 216 boring bin. This confirms that the largest group splits when the scale is not
 217 perfect. When the scale differs significantly from the best value, either to the
 218 left of right, the solid curve drops to two matches¹ and the dotted curve shows
 219 only one match.

220 Given the best scale, we use the corresponding 2D histogram to find the
 221 matches aggregated in the largest group at this scale. Figure 8 shows the
 222 correct and incorrect matches.

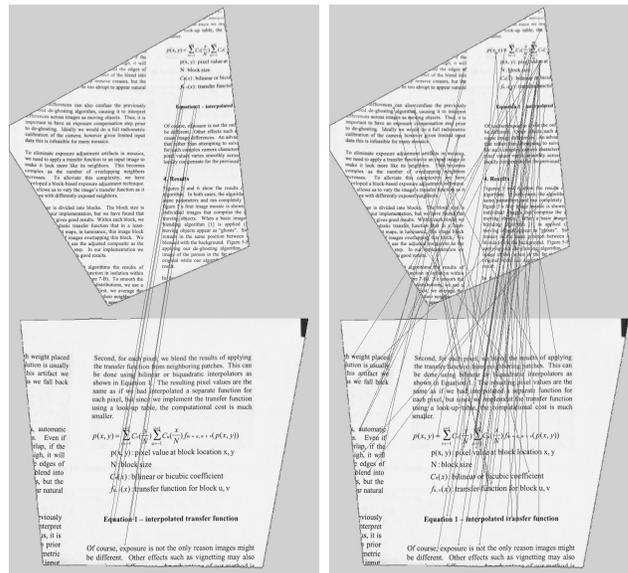


Fig. 8. Correct (left) and incorrect (right) matches in the PCA-SIFT result.

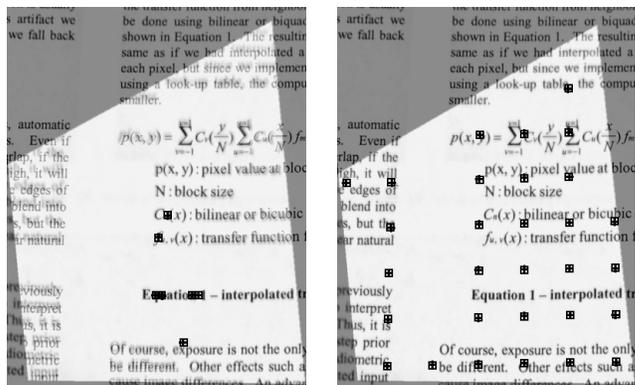
¹ The largest group has two matches because one pair of matched points is duplicated in the output of PCA-SIFT.

223 In the above analysis, the concept of compatibility groups is similar to the
224 compatibility of matches in [3]. In a broader view, we can see that *voting*
225 lies at the heart of our method, [3] and other RANSAC variations, and many
226 geometric hashing based methods. To deal with outliers, RANSAC relies on
227 the “random” chance of picking a good set, which could be inefficient or in-
228 effective as the chance decreases when outliers dominate; geometric hashing
229 solves the efficiency problem by doing the majority of computation off-line us-
230 ing training data (entire images, not local features). In the area of document
231 mosaicing, our method takes advantage of the fact that image rectification
232 can first remove a great part of the uncertainties, and as a result the voting
233 itself becomes deterministic and efficient.

234 Taking the correct matches, we compute an initial projective transformation
235 between the two images and map one into the other, as shown in Figure 9(a).
236 However, because good matches tend to reside near the overlapped region’s
237 center, the registration is inaccurate near the border. We further refine the
238 registration using cross-correlation block matching. This results in a dense
239 and accurate set of matched points covering the whole overlapped area, which
240 allows us to compute a refined projective transformation (see Figure 9(b)).

241 4 Seamless Composition

242 As we have stated in the introduction, there are three difficulties in creating
243 a seamless document mosaic. The first is due to inconsistent lighting across
244 two images. Conventional blending does not address overall lighting incon-
245 sistency, and it works well for general photos only because people accept
246 lighting changes in natural scenes. However, documents are fundamentally



(a)

(b)

Fig. 9. Image registration results where squares and crosses indicate the matched points in two images. (a) Registration using correct PCA-SIFT matches shows misalignment near borders of overlapping region. (b) Registration using additional matches obtained from cross-correlation block matching is very accurate.

247 binary with black print on white paper, and viewers' eyes are very sensitive
 248 to varying shade in documents. Typically, the histogram of a document im-
 249 age is bimodal. Different lighting conditions cause the two modes to shift.
 250 One way of balancing the lighting across two document images is to binarize
 251 both images. However, binarization introduces artifacts. Instead, we choose
 252 localized histogram normalization. The basic idea is to compute the local his-
 253 togram in a small neighborhood, normalize the histogram such that the two
 254 modes are transformed to black and white respectively (or very dark and
 255 light gray). Histogram normalization preserves the transition between back-
 256 ground and foreground, so the result is more pleasing to view. For documents
 257 containing grayscale or color contents, one choice would be to apply segmen-
 258 tation first, then compute histogram normalization parameters in the bimodal
 259 areas and estimate these parameters in grayscale or color areas via interpola-
 260 tion/expolation.

261 The second problem is registration error, and the third is uneven sharpness of

262 patch images. We solve both with *selective image composition*, i.e., each pixel
263 in the result is chosen from the image with the best sharpness. We measure
264 sharpness in an image by the local average of gradient magnitudes. In the
265 following, the index of the selected image for a pixel is called the *decision* for
266 this pixel.

267 The pixel-level decisions can be represented by a map in which the same
268 decisions are grouped into regions. The boundaries of decision regions may
269 intersect characters and words. Thus, if we apply pixel-level decisions directly,
270 some characters or words may consist of pieces with different sharpness chosen
271 from different images, which is not desirable. Furthermore, mis-registration
272 tends to break decision regions into small pieces, resulting in ‘ghost’ images.

273 Therefore we aggregate the pixel-level decisions at the word level. This re-
274 quires finding words. To do this, we compose an averaging image for the over-
275 lapped area, binarize it, dilate the foreground, and find connected components.
276 The dilation has two effects. First, areas that may contain ‘ghost’ images are
277 merged into the nearest component. Second, the width of our dilation kernel
278 is set to be larger than its height, so components of a word are more likely to
279 be merged than components from upper or lower text lines. As a result, most
280 connected components contain a word. Next, all the pixels inside a connected
281 component vote with their pixel-level decision, and the majority vote is taken
282 as the component decision. The values for all the pixels of the component are
283 selected from the winning image. This process ensures that ‘ghost’ images are
284 eliminated and words do not have an uneven sharpness. For background areas
285 (areas that are not included in word regions), the variation of sharpness is not
286 visible, so we use the pixel-level decisions directly (without voting) to assign
287 their values.

288 Figure 10 illustrates the process of selective image composition and the re-
289 sults. Figure 10(a) shows that most components consist of a single word. Fig-
290 ure 10(b) shows the component-level decision map by two shades of gray. The
291 arrows indicate words that are cut into different parts in pixel-level decision
292 map but not broken in component-level decision map.

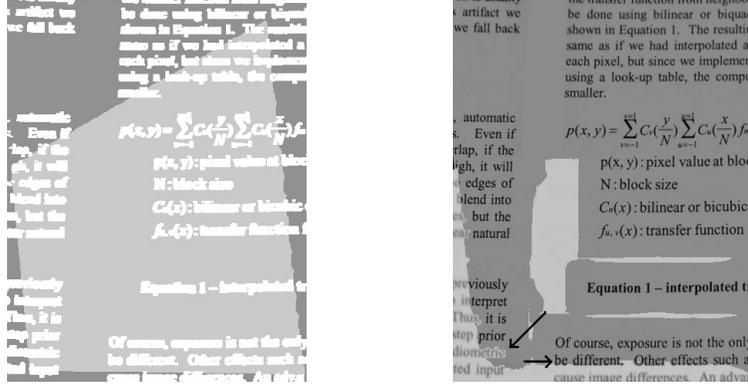
293 Words may still be broken by the boundaries of the overlapping area (e.g.,
294 “previously” and “interpret” in the lower left part of Figure 10(b)). In this
295 case, half of the broken word has better sharpness than the other half. One
296 could select the entire word from the image with lower sharpness to eliminate
297 the difference in sharpness. This choice depends on user preference.

298 In the background area, the pixel-level decisions result in a large light gray
299 region embedded in dark gray area. This does not create visible differences in
300 the final image because the variation of sharpness in the background is small.
301 In Figure 10, the comparison between (c) and (d) shows that our approach
302 preserves the sharpness. In Figure 10(e), the overlapping area boundary is
303 visible. It is eliminated in Figure 10(f).

304 5 Experiments

305 A quantitative evaluation of the rectification step (using synthetic data) is
306 given in [6]. In this paper, we evaluate overall results on real images.

307 For each test document, we obtained a scan at 300 dpi and use the scan as
308 the ground truth image. We took pictures of the document with a Canon
309 EOS 300D digital camera (6M pixels) and a 28-80mm (35mm equivalence)
310 zoom lens. All camera settings were put in auto-mode. First, we posted the



(a)

(b)

$$\sum_{v=-1}^{v=1} C_v\left(\frac{y}{N}\right) \sum_{u=-1}^{u=1} C_u\left(\frac{x}{N}\right) f_{v,u}(x,y)$$

(c)

$$\sum_{v=-1}^{v=1} C_v\left(\frac{y}{N}\right) \sum_{u=-1}^{u=1} C_u\left(\frac{x}{N}\right) f_{v,u}(x,y)$$

(d)

ation 1. The result is the same as if we had interpolated at each pixel, but since we implemented using a look-up table, the computation is smaller.

(e)

(f)

Fig. 10. Selective image composition. (a) Connected components are represented by white. The overlapping area is represented by light gray. (b) The binary selection decision map distinguished by dark and light gray. (c, e) Weighted averaging result. (d, f) Selective image composition result.

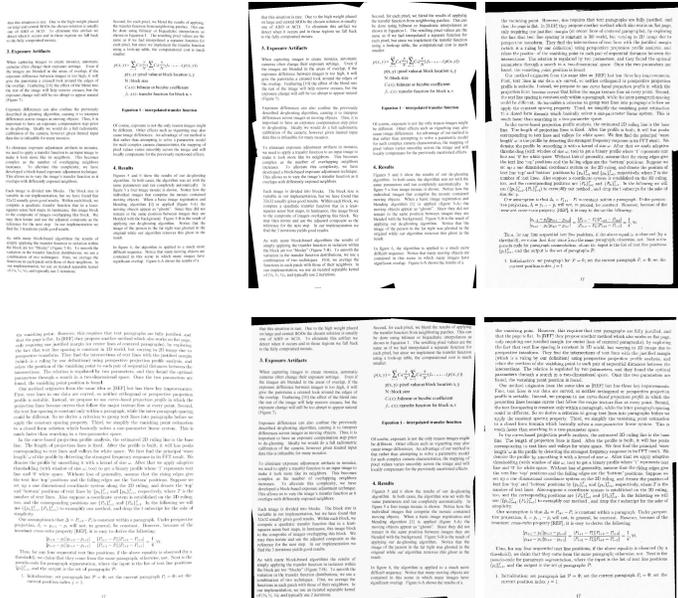
311 document on a wall, put the camera on a tripod, and carefully calibrated it
 312 so that the image is perspective free. We provided sufficient ambient light so
 313 that a small aperture (f5.6) could be used without flash. The result image
 314 represents, in some sense, the optimal non-mosaic image we can expect from
 315 the camera. Then we lowered the ambient light to typical indoor level (which
 316 caused the camera flash to go off in many cases), placed the document on
 317 a desk, sat at the desk, held the camera in hand and took a set of pictures
 318 with various angles and zooms. These image patches are fed to our mosaicing
 319 algorithm. Figure 1 contains three example image patches.

320 In the first round, we collected four documents from scientific journals and
 321 conference proceedings. We captured four patches for each document, making

322 sure that they all overlap each other. The perspective in the patches are kept
323 moderate to minimum. This represents the scenario where the user has some
324 control over the pose but needs higher resolution than a global view. We fed the
325 four patches in all possible orders ($C_4^4 = 24$) to our mosaicing module. During
326 the process, the first image serves as the initial composite, and each patch is
327 registered to the current composite which contains all previous patches.

328 In the second round, we captured the images with less constraints to evaluate
329 the limits of applicability of our approach. The number of patches varies from
330 eight to twelve. Compared to the first round, the patches cover smaller areas,
331 and perspective is from moderate to high. This simulates the case where a
332 low resolution camera is used and it is difficult to control the document pose.
333 We visually inspected and recorded which patches overlap each other. Then
334 ten random orderings of the patches were generated under the condition that
335 each patch overlaps with at least one patch before it. This ensures that when a
336 patch is fed to the mosaicing module, it overlaps with the composite generated
337 so far.

338 We find no significant difference between the result in the first round and the
339 second round, except for expected differences of resolution. Figure 11 shows
340 the camera-captured frontal views of two documents, and two composite im-
341 ages for each document. The global views in Figure 11(a) are enhanced by
342 the same local histogram normalization used in document mosaicing. Slight
343 barrel effects are visible, due to the wide lens distortion. All composites in
344 Figure 11(b) have higher resolution and do not show barrel effect, compared
345 to (a). The enlarged views show two portions near the border of underlying
346 patches. There are some fuzziness, misalignment, and “ghost” effects. Never-
347 theless, the border itself is undetectable.



(a)

(b)

(c)

Fig. 11. Perspective-free images compared to composite images. (a) Full frontal views captured by a camera. (b) Two composite images for each document. (c) Enlarged portions of composites near borders underlying patches.

348 Except for visual inspection which could be subjective, we also compared the
 349 mosaics to the global views quantitatively in an overall sense. Since our sam-
 350 ples are mostly text documents, we used OCR as the image quality appraiser.
 351 For each document, we applied OCR to the digital scan, the camera-captured
 352 global view, and all the composites. We used the OCR text from the scan as
 353 the ground truth, against which we computed the character and word recogni-
 354 tion rates for the global views and composites, respectively. Table 5 shows the
 355 number averaged over all documents. The OCR performance on composites is
 356 very close to that on the perspective free global view.

357 The PCA-SIFT used in our experiment is trained using generic image data.
 358 We also tried document images as training data. However, no difference was
 359 found in their performance in percentage of false alarms.

meter space.
 analysis, the estimat
 d. After the prof
 white space. We
 ngest frequency re
 kernel of size ω .

stant within a par
 al, be constant.
 to derive the follo

$$\frac{-P_i || P_{i+3} - P_{i+2} |}{-P_i || P_{i+3} - P_{i+1} |}$$

	Global views	Composite images
Character recognition rate	92.3%	91.0%
Word recognition rate	89.2%	89.5%

Table 1

Average OCR rates of global views and composite views.

360 In our experiments, computing a mosaic could be very time consuming. De-
361 pending on hardware and image size, it may take up to ten minutes. This is
362 partially because our prototype is built in MATLAB and not optimized for
363 speed. The most demanding part is rectification, especially the texture flow
364 computing. The blending part comes second. PCA-SIFT comes third. The reg-
365 istration step is negligible. Overall, the complexity is roughly linear in terms
366 the numbers of pixels.

367 6 Summary

368 In this paper we demonstrate a document mosaicing method which deals with
369 severe perspective distortion, large displacement and small overlapping area.
370 The first step, geometric rectification, greatly reduces the complexity of the
371 registration problem. The second step, registration, is robust against large
372 number of outliers found by feature point matching algorithms. The last step,
373 blending, composes a seamless, “ghost” free mosaic with optimal sharpness.
374 While the rectification step only works on text areas in documents, the other
375 two steps can be applied to non-text images without significant modifications.

376 **References**

- 377 [1] P. J. Burt and E. H. Adelson. A multiresolution spline with application to
378 image mosaics. *ACM Transactions on Graphics*, 2(4):217–236, 1983.
- 379 [2] S. Gold and A. Rangarajan. A graduated assignment algorithm for graph
380 matching. *IEEE Trans. PAMI*, 18(4):377–388, April 1996.
- 381 [3] F. Isgrò and M. Pilu. A fast and robust image registration method based on an
382 early consensus paradigm. *Pattern Recognition Letters*, 25(8):943–954, 2004.
- 383 [4] Y. Ke and R. Sukthankar. PCA-SIFT: a more distinctive representation for
384 local image descriptors. In *Proc. CVPR*, volume 2, pages 506–513, 2004.
- 385 [5] J. Liang, D. DeMenthon, and D. Doermann. Camera-based document image
386 mosaicing. In *Proc. ICPR*, 2006.
- 387 [6] J. Liang, D. DeMenthon, and D. Doermann. Geometric rectification of camera-
388 captured document images. *IEEE Trans. PAMI*, 30(4):291–605, Apr. 2008.
- 389 [7] M. Mirmehdi, P. Clark, and J. Lam. Extracting low resolution text with an
390 active camera for OCR. In *Proc. IX Spanish Sym. Pat. Rec. and Image Proc.*,
391 pages 43–48, May 2001.
- 392 [8] T. Nakao, A. Kashitani, and A. Kaneyoshi. Scanning a document with a small
393 camera attached to a mouse. In *Proc. WACV'98*, pages 63–68, 1998.
- 394 [9] B. S. Reddy and B. N. Chatterji. An FFT-based technique for translation,
395 rotation, and scale-invariant image registration. *IEEE Trans. Image Proc.*,
396 5(8):1266–1271, 1996.
- 397 [10] T. Sato, A. Iketani, S. Ikeda, M. Kanbara, N. Nakajima, and N. Yokoya. Mobile
398 video mosaicing system for flat and curved documents. In *Proceedings of 1st*
399 *International Workshop on Mobile Vision*, pages 78–92, 2006.

- 400 [11] T. Sato, A. Iketani, S. Ikeda, M. Kanbara, N. Nakajima, and N. Yokoya. Video
401 mosaicing for curved documents based on structure from motion. In *Proc.*
402 *ICPR*, volume 4, pages 391–396, 2006.
- 403 [12] K. Schutte and A. M. Vossepoel. Accurate mosaicking of scanned maps, or how
404 to generate a virtual A0 scanner. In *Proc. ASCI'95*, pages 353–359, 1995.
- 405 [13] M. I. Tomohiro Nakai, Koichi Kise. Use of affine invariants in locally
406 likely arrangement hashing for camera-based document image retrieval. In
407 *International Workshop on Document Analysis Systems*, pages 541–552, 2006.
- 408 [14] A. P. Whichello and H. Yan. Document image mosaicing. In *Proc. ICPR*, pages
409 1081–1083, 1998.
- 410 [15] A. Zappala, A. Gee, and M. J. Taylor. Document mosaicing. *Image and Vision*
411 *Computing*, 17(8):585–595, 1999.