# PARTIAL CONTOUR MATCHING FOR DOCUMENT PIECES WITH CONTENT-BASED PRIOR

*Fabian Richter, Christian X. Ries, Stefan Romberg, Rainer Lienhart*

Multimedia Computing and Computer Vision Lab, University of Augsburg, Germany
{richter, ries, romberg, lienhart}@informatik.uni-augsburg.de

## ABSTRACT

In this paper we present a method for aligning shredded document pieces based on outer contours and content-based prior information. Our approach relies on domain-specific knowledge that document pieces must complement each other when aligned correctly. Building on this intuition we propose a variant of MSAC (M-estimator SAmple Consensus) to estimate an hypothesis that recovers the spatial relationship between pairs of pieces. To do so we first approximate their boundaries by polygons from which we define consensus sets between fragments. Each consensus set provides multiple hypotheses for aligning one piece onto the other. An optimal hypothesis is identified by applying a two-stage procedure in which we discard locally inconsistent hypotheses before verifying the remainder for global consistency.

*Index Terms*— Document analysis, MSAC, partial contour matching

## 1. INTRODUCTION

There are many reasons for making sensitive data in documents unreadable including tax fraud, business crime, or other criminal intentions. Our work thus investigates on recovering hand-torn document pages, for which we propose a method to precisely align pairs of pieces based on their boundary and content-based prior information. A system capable of recovering the relative layout of document pieces could be of great use in many real-world applications. For instance, a recent project of the Fraunhofer Institute [1] deals with the problem of reassembling documents related to the Stasi, which was the former secret police of the GDR. Shortly before the end of the Socialist regime of the GDR in 1989, members of the Stasi destroyed millions of documents containing evidence about their activities. Many of those files were simply torn by hand and thus can be attempted to be recovered. Another example is that of an Israeli group who recently scanned and restored 250,000 historic documents that are hundreds of years old.

For this reason we propose a method that could guide human users in the reassembling process. Common requirements of such a system include (i) the ability to handle rotated pieces and (ii) robustness to image noise stemming from the



**Fig. 1**. Example for polygons (grey line segments) that approximate the pieces' outer contours. Pairs of support points among the fragments' boundaries form a consensus set (green dots). This set tends to contain many correct one-to-one correspondences (inliers) if the contour segments match. Based on pairs of inliers, e.g. $(i, j)$ and $(i', j')$, we may then determine an hypothesis that precisely aligns the two fragments.

digitizing process. Our goal is to recover the spatial relationship for pairs of arbitrarily oriented fragments from a dataset of hand-shredded documents. To obtain such partial solutions we have to estimate the orientation and relative position of pairs of pieces, without knowing the intact page in advance.

To accomplish this, we approximate the boundary of each fragment by a simple polygon to obtain a less complex representation. In this respect, identifying partially matching contour segments from pairs of document pieces comes down to finding hypotheses that satisfy geometric constraints based on the pieces' outer contours. As illustrated in figure 1, the basic idea is to construct consensus sets which define one-to-one correspondences between support points of the polygons. We then use these correspondences to determine hypotheses that embed boundaries of pieces into a common coordinate system in a way that they complement each other. Most importantly, since content overlap provides evidence for hypotheses being incorrect, we capitalize on this prior knowledge to identify and discard invalid transformations at early stages of the alignment process.

## 2. RELATED WORK

A popular approach to partial contour matching is the Smith-Waterman algorithm [2], which is a dynamic programming technique initially proposed to find maximally homologous subsequences among two molecular sequences. Lately many approaches also make use of this algorithm for shape recognition. For instance, Bunke and Bühler [3] devised a method for recognizing arbitrary two-dimensional shapes which is invariant under translation, rotation and partial occlusion. A similar route was taken by Chen et al. [4], who also apply the Smith-Waterman algorithm to evaluate the similarity of two shapes. Their work differs from others in that they use a probabilistic similarity measure instead of distance functions.

One approach similar to our work is that of Stieber et al. [5], who adopt the Smith-Waterman algorithm for partial contour matching between document pieces. In their method the authors only compare contour segments delimited by subsequent corners, which helps in avoiding false matches that are inconsistent regarding the pieces' global geometry. Similar as in our work, they also incorporate geometric information into their similarity score to assess the quality of an individual alignment. We note that our approach does not require corner detection, which could be beneficial for cases where the pieces' corners can not be identified without ambiguity.

Another related approach has been proposed by Donoser et al. [6], who sample points from object silhouettes to represent shapes. Their method utilizes local and global geometry for detecting subparts of two shapes that possess high similarity. Since their matching procedure is formulated as order-preserving assignment problem, it is conceptually similar as to how we form consensus sets from sets of support points.

Other recent works investigate on the automatic assembly of documents [7] and photos [8]. The latter approach of Cao et al. deals with the problem of simultaneously reconstructing multiple photos from a collection of pieces. The authors use a curve matching algorithm to first identify matching pairs of pieces, which are then clustered into groups of pieces from a single photo. Finally, to assemble pieces into an intact photo, a spanning tree algorithm is applied that also verifies the validity of the solution in terms of geometric consistency.

## 3. PRELIMINARIES

We commence by introducing the dataset used in this work before describing our approach in section 4. We then explain in section 5 how content-based information about the pieces' orientation is used to discard incorrect hypotheses early in the alignment process.

### 3.1. Dataset of Shredded Documents

In this work we use the *bdw082010* dataset [7], which consists of 96 document pages that have been shredded by hand into 16 pieces each. Our preprocessing closely follows [7] to obtain an approximation of each fragment's contour. Since each fragment comes with a binary segmentation mask that identifies its foreground, we first determine the set of all boundary points $P$ using the algorithm of Suzuki et al. [9]. Afterwards we apply the Douglas-Peucker algorithm [10] to find a small subset of *support points* $\hat{P} = \{\hat{\boldsymbol{p}}_1, \ldots, \hat{\boldsymbol{p}}_n\} \subseteq P$ that constitutes a less complex description of each piece's outer contour. By connecting consecutive pairs of support points with line segments one finally obtains a polygon which approximates the exact contour up to a predefined precision.

For the remainder of this work we write $s \equiv (P^s, \hat{P}^s)$ in short for pieces having their support points ordered in clockwise direction, while those being processed counterclockwise are denoted by $t$. Based on this notion of ordered boundary points our goal is to identify those *correspondences* between two fragments $s$ and $t$ that originate from the same location in the former document. Throughout the following sections we will denote a correspondence by $(i, j)$, which associates the $i$-th support point from the contour of piece $s$ with the $j$-th point on the contour of $t$.

## 4. APPROACH

We now describe our approach to find and align contour segments that are likely to adjoin each other in an intact document. Our method *2p–PCM–MSAC* is a variant of MSAC [11] (M-estimator SAmple Consensus), which estimates parameters of a model (hypothesis) from a set of observations.

In our problem setting, any correspondence between the contours of two fragments is considered an observation. Naturally, only very few correspondences stem from the same location in the document (*inliers*), while the majority of boundary points is spatially disconnected in the document (*outliers*). For partial contour matching (*PCM*), we thus aim to find the model that best aligns the inliers onto each other by first performing a translation, followed by an in-plane rotation. Since this group of Euclidean transformations has 3 degrees of freedom, it can be estimated from only two pairs of points – hence the prefix *2p-PCM*.

**Futility of Unbiased Random Sampling.** To estimate a correct model we thus need to find one pair of inliers from the set of all observations. In RANSAC [12] (RANdom SAmple and Consensus) this problem is approached by iteratively sampling a random subset of data points that could possibly be inliers. Given a desired probability $p$ that RANSAC finds at least $n = 2$ inliers, the number of iterations that are needed is upper bounded by $k = \log{(1 - p)} / \log{(1 - w^n)}$. However, due to the inherently small inlier ratio $w$, we note that random sampling is not the method of choice in our scenario. Even for conservatively estimated parameters, the required number of iterations would be in the order of $10^6$.

## 4.1. Building the Consensus Set

Instead of random sampling we now propose a more effective method to find correspondences that are likely inliers. Throughout algorithm 1 (line 4) we once consider each correspondence between two fragments $s$ and $t$. Hence, among all combinations of support points, we encounter at least one point pair that corresponds to an inlier. As exemplified in figure 1, pair $(i, j)$ is an inlier, because both points refer to the same position in the document. With regard to this initial correspondence, called a pair of *anchor points*, any second pair $(i', j')$ can also be considered an inlier if the line segments enclosed by $i$ and $i'$ on $s$ and $j$ and $j'$ on $t$ are similar.

For partial contour matching there are two notions of similarity that come into mind. One necessary condition for line segments to match is that they have equal *length*. On the other hand, both segments must also be similar in terms of *shape*. We address the latter aspect in section 4.4, where the validity of each hypothesis is verified in terms of local geometry.

In the following we first focus on finding equally long line segments among the two fragments, from which we then form *consensus sets*. For this reason we precompute the absolute distances of points on the fragments' contours with respect to fixed reference points. Hence all relative distances between two arbitrary points on the same contour can be computed on demand by performing two lookups. In our algorithm (line 5), the relative distances w.r.t. anchor points $i$ and $j$ are indexed in $\tau_s$ and $\tau_t$, respectively.

Based on the precomputed distances, the idea is to search for points on both sides that have approximately the same distance with respect to their anchor points. To accomplish this we only consider short *contour segments*, in clockwise direction on $s$ and counterclockwise direction on $t$, respectively. Each of these segments is restricted to cover at most 20% of the two fragments' total boundary length, as illustrated in figure 1 (marked regions). From all possible point combinations between those two segments we then greedily choose a subset of equidistant one-to-one correspondences, which are indicated in the figure by identically colored dots. Correspondences having dissimilar euclidean distances to their respective anchor points are discarded during this process, as this typically indicates an outlier.

All established correspondences $(i', j')$ are candidates for inliers and thus become part of consensus set $C_{ij}$ (line 6). An important observation is that the cardinality of this set always depends on the anchor points. For instance, when correspondence $(i, j)$ itself is no inlier, its consensus set tends to contain only very few (if any) equidistant pairs of points.

## 4.2. Hypotheses from Consensus Sets

We now describe how to obtain multiple hypotheses (line 9) from each consensus set, which are later checked for inconsistencies regarding local geometry (see section 4.4).

---

### Algorithm 1: 2p–PCM–MSAC

**Input** : $s \equiv (P^s, \hat{P}^s), t \equiv (P^t, \hat{P}^t)$
**Output**: model $h_2^*$

**Initialization**
2    $\mathcal{H}^1 = \emptyset, \mathcal{H}^2 = \emptyset$

**Step1 (Creating locally verified hypotheses)**
4    **foreach** $(i, j) \leftarrow \hat{P}^s \times \hat{P}^t$ **do**
5      $[\tau_s, \tau_t] \leftarrow$ `rel-distances` $(i, j)$
6      $C_{ij} \leftarrow$ `consensus-set` $(i, j, \tau_s, \tau_t)$
7      $\mathcal{H}_{ij} = \emptyset$
8      **foreach** $c_k \in C_{ij}$ **do**
9        $h \equiv (\boldsymbol{t}, \alpha_k) \leftarrow$ `compute-h` $(i, j, c_k)$
10        **if** `valid-geometry` $(h, C_{ij})$ **then**
11          $\mathcal{H}_{ij} = \mathcal{H}_{ij} \cup \{h\}$
12      $h_1^* \leftarrow \operatorname{argmin}_{h \in \mathcal{H}_{ij}} \{\epsilon_1(h; C_{ij})\}$
13      $\mathcal{H}^1 = \mathcal{H}^1 \cup \{h_1^*\}$

**Step2 (Global verification)**
15    **foreach** $h \in \mathcal{H}^1$ **do**
16      **if** `valid-geometry` $(h, (\hat{P}^s, \hat{P}^t))$ **then**
17        $\mathcal{H}^2 = \mathcal{H}^2 \cup \{h\}$

18    $h_2^* \leftarrow \operatorname{argmin}_{h \in \mathcal{H}^2} \{\epsilon_2(h)\}$
   **return** $h_2^*$

---

Given the coordinates of two anchor points, we first compute an offset vector $\boldsymbol{t}$ that translates support point $j$ onto $i$ in a common coordinate system. Since both points coincide in $i$ after translation, we also use this point as rotation center for all hypotheses. From each pair of points $c_k \equiv (i', j')$ in consensus set $C_{ij}$ we then determine a rotation angle $\alpha_k$, based on the line segments between $i, i'$ and $j, j'$, respectively. This fully defines one Euclidean transformation $h \equiv (\boldsymbol{t}, \alpha_k)$ for each member in $C_{ij}$.

### 4.3. Signed Nearest Distance

As will become clear in the following section, we further have to evaluate distances of points to a fragment's contour. Given the set of all support points $\hat{P}$ that define the vertices of the boundary, we measure the minimal perpendicular distance of point $\boldsymbol{p}$ to any of the polygon's line segments. We denote this distance function by $d(\boldsymbol{p}; \hat{P})$. Finally, the sign of the distance is determined, that is whether the point lies inside (positive) or outside the polygon (negative).

### 4.4. Local Geometric Verification

At this point we want to utilize local geometry to perform a fast spatial verification of each hypothesis found thus far.
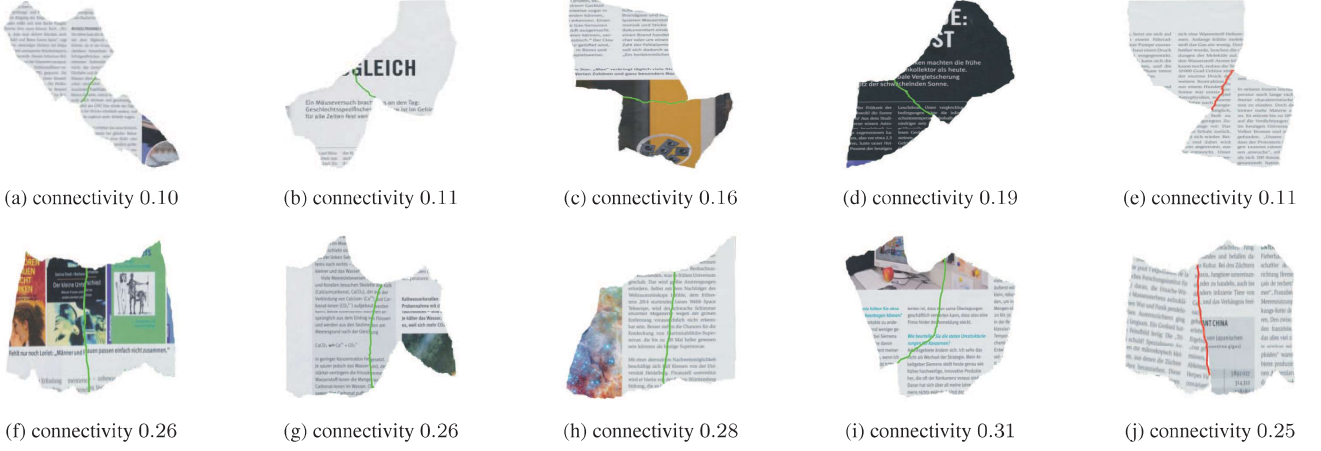
(a) connectivity 0.10    (b) connectivity 0.11    (c) connectivity 0.16    (d) connectivity 0.19    (e) connectivity 0.11

(f) connectivity 0.26    (g) connectivity 0.26    (h) connectivity 0.28    (i) connectivity 0.31    (j) connectivity 0.25

**Fig. 2**. Results for different degrees of connectivity. **Top row:** Rectified examples for correctly aligned pieces (min-overlap greater than 0.5), with *low–medium* connectivity that increases from left to right (a)-(d). **Bottom row:** Results for examples with *medium–high* connectivity (see (f)-(i)). Examples for incorrect solutions, e.g. pairs having zero min-overlap, are shown in the rightmost columns (e) and (j).

Based on our conception that correct models never produce any overlapping content between aligned fragments, we discard hypotheses that map any point from their consensus set into the foreground region of the piece's counterpart. To determine whether this is the case we use the signed distance of each point in $C_{ij}$ to the nearest boundary point of the other fragment. Let the coordinates for a pair of points $c_k \equiv (i', j')$ be denoted by $\hat{P}^s[i']$ and $\hat{P}^t[j']$. Then, if either

$$\max_{c_k \in C_{ij}} \left\{ d\big(\hat{P}^s[i'], h(\hat{P}^t)\big) \right\} > T \quad (1)$$

or

$$\max_{c_k \in C_{ij}} \left\{ d\big(h(\hat{P}^t[j']), \hat{P}^s\big) \right\} > T \quad (2)$$

exceeds threshold $T$ we do not add $h$ to $\mathcal{H}_{ij}$ (see line 10 in alg. 1). Note that $h(\cdot)$ refers to transformed coordinates after applying $h$. All remaining hypotheses are then ranked according to their accuracy in aligning points from the consensus set. For this purpose we compute a (truncated) squared euclidean distance between transformed points from $t$ and their counterparts from $s$ as follows:

$$\epsilon_1(h; C_{ij}) = \sum_{c_k \in C_{ij}} \left\lfloor \|\hat{P}^s[i'] - h(\hat{P}^t[j'])\|^2, D_{max} \right\rfloor \quad (3)$$

We write $\lfloor \cdot, \cdot \rfloor$ and $\lceil \cdot, \cdot \rceil$ in short for the min- and max-operation, respectively. From all initial models we finally add only the best scoring $h_1^* = \operatorname{argmin}_{h \in \mathcal{H}_{ij}} \{\epsilon_1(h; C_{ij})\}$ to set $\mathcal{H}^1$ of locally verified hypotheses.

### 4.5. Global Geometric Verification

We recap that up to this point, only correspondences within consensus sets have been used to preemptively discard wrong alignment models. After this initial verification, we must now validate only very few hypotheses in terms of global geometry, where pieces are dealt with as a whole. For this purpose, we first use the same criteria as for our local verification step (eq. (1), (2)), however, we now check all support points from $s$ and $t$ for content-overlap. Only those hypotheses that also pass this global verification (line 16 in alg. 1) are retained in set $\mathcal{H}^2$. Finally, our aim is to select only one hypothesis from this set that best aligns the two fragments onto each other.

To accomplish this, all remaining models must be ranked according to our conception of an optimal alignment result. Hence we introduce the following function that, individually for pieces $s$ and $t$, penalizes content overlap as well as spatially disconnected boundaries:

$$\epsilon_2(h) = \sum_{i=1}^{n_s} \delta\big(\hat{P}^s[i], h(P^t)\big) + \sum_{j=1}^{n_t} \delta\big(h(\hat{P}^t[j]), P^s\big) \quad (4)$$

Variables $n_s$ and $n_t$ correspond to the number of support points on the two fragments. Using $d \equiv d(\hat{p}, P)$ as substitute for the signed nearest distance (see section 4.3) we define:

$$\delta(\hat{p}, P) = \left\lceil \lfloor d, D_o \rfloor^{e_o}, 0 \right\rceil + \left\lceil \lfloor -d, D_b \rfloor^{e_b}, 0 \right\rceil \quad (5)$$

Considering all remaining hypotheses we finally choose $h_2^* = \operatorname{argmin}_{h \in \mathcal{H}^2} \{\epsilon_2(h)\}$ as the best model to align fragment $t$ onto $s$. Parameters $D_o$ and $D_b$ control the maximum distance of each point to the boundary before it is considered as outlier. Besides, exponents $e_o$ and $e_b$ control the behavior of the error function regarding points that lie on the inside and outside of the polygons, respectively. Either of the two criteria can be emphasized by using different exponents, in which case the error function becomes asymmetric. Note that all parameters were chosen empirically, based on qualitative results on the training and validation set.
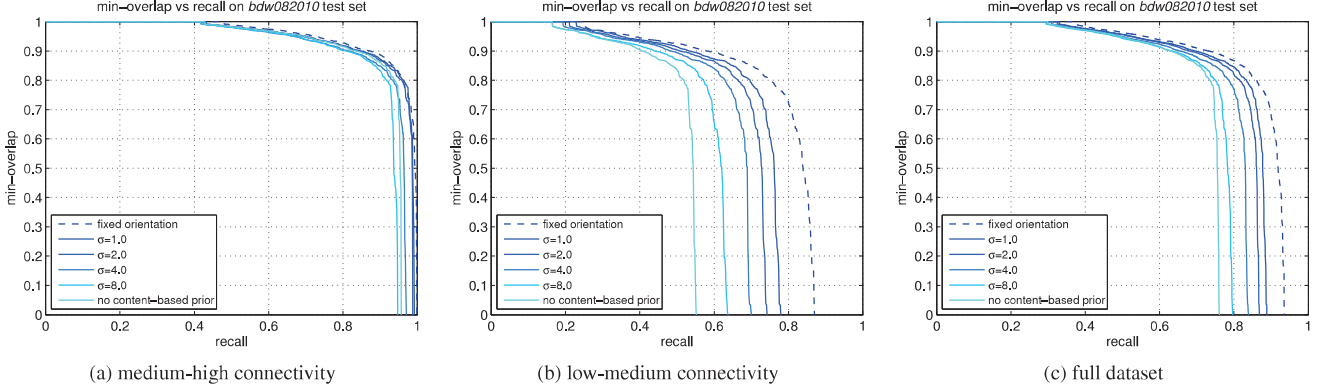
(a) medium-high connectivity      (b) low-medium connectivity      (c) full dataset

**Fig. 3**. Performance evaluation in terms of min-overlap vs. recall, for different levels of connectivity (figure 3a and 3b), as well as for the full dataset (figure 3c).

## 5. CONTENT-BASED PRIOR INFORMATION

For text pages, content-based prior information about the orientation of documents pieces can be obtained by running a text detector (e.g. [13]) in advance. In the following we assume that despite the absolute orientations $\theta_s$ and $\theta_t$ for pieces $s$ and $t$ are unknown, one may still infer an estimate for their relative orientation from the detector output. By running the detector individually for both pieces, one obtains an estimate for the true orientations and their enclosed angle $\hat{\alpha}$.

In our algorithm we use this estimate to discard all hypotheses that do not conform with the predicted relative orientation. For this purpose we model the uncertainty associated with a detector as Gaussian noise with standard deviation $\sigma$. Therefore, any hypothesis that predicts a relative orientation $\alpha_k \notin [\hat{\alpha} - K\sigma, \hat{\alpha} + K\sigma]$ is omitted. We note that estimate $\alpha_k$ is presumably only correct up to $180°$, because a detector is likely to output the text orientation, but not its direction. For our experiments in section 6.1 we set $K$ empirically to balance the number of hypotheses that can be discarded without falsely omitting correct hypotheses.

## 6. EVALUATION

In this section we discuss the performance criteria used for our experiments on the bdw082010 dataset, which has been described in section 3.1. For our evaluation we first introduce the notion of *connectivity*, which can be considered a measure for the adjacency of two pieces in the intact document.

**Connectivity.** In this work the connectivity of two pieces is defined as the length of the adjacent boundary segments relative to their overall contour length. Intuitively, pieces having a high connectivity provide significant evidence for being adjacent in the document. As a consequence, finding the correct contour segments for examples with low connectivity is inherently more difficult, as is reflected by our first experiment in

section 6.1. Some qualitative results for examples with different degree of connectivity are depicted in figure 2. We would like to point out that all of those fragments have been aligned according to the hypotheses resulting from algorithm 1.

**Overlap between Contour Segments.** For any predicted hypothesis that aligns two pieces $s$ and $t$, we infer contour segments $\hat{l}_s$ and $\hat{l}_t$, one along each fragment's boundary, where pieces complement each other (e.g. green and red lines in figure 2). Since correct line segments $l_s$ and $l_t$ are known from ground truth data, we define their pairwise overlap as the minimum intersection over union from both sides, i.e.

$$min\text{-}overlap(\hat{l}_s, \hat{l}_t) = \left\lfloor \frac{\hat{l}_s \cap l_s}{\hat{l}_s \cup l_s}, \frac{\hat{l}_t \cap l_t}{\hat{l}_t \cup l_t} \right\rfloor. \qquad (6)$$

The min-overlap evaluates to 1 for two predicted line segments only if they are a perfect match regarding ground truth data. This measure is suitable for assessing the quality of an hypothesis, because it gives gradually smaller values if either of the segments is not located correctly, while likewise penalizing segments that are too long or short.

### 6.1. Experiments

Our first experiment evaluates the effectiveness of the proposed method in terms of min-overlap vs. recall. We compute an upper bound for the maximal achievable performance by assuming that the absolute orientation of fragments is known and held fixed, so that each hypothesis estimates only a translation. In this scenario we achieve a recall of $93.6\%$ on the full dataset, as can be seen in figure 3c. With higher degree of connectivity, the recall increases from $87.1\%$ for hard examples (figure 3b) to $99.9\%$ for relatively ones (figure 3a).

In a second scenario we assume that the orientation of pieces is unknown, but can be estimated from content-based information in advance (see section 5). Because a text detector only allows for an approximately correct orientation as-
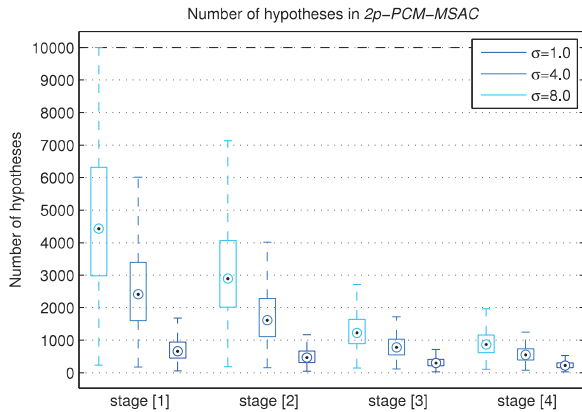
**Fig. 4**. Number of hypotheses at different stages of 2p–PCM–MSAC, for varying $\sigma$. Each group consists of 3 boxplots with their median value (circle) as well as the 25th and 75th percentiles. See text for details on stages 1–4.

signment, we plot overlap-recall curves for different levels of uncertainty $\sigma \in \{1, 2, 4, 8\}$ (in degrees) about this estimate. As can be seen in figures 3a and 3b, having a more accurate prediction for the pieces' orientation (smaller $\sigma$) constantly improves the performance in all cases. However, this is more pronounced for examples with low connectivity. Finally we also run an experiment without using any content-based prior information, in which case we still achieve a recall of 76.1% on the full dataset (see figure 3c). Examples for aligned fragments are given in figure 2. Pieces with less than four inliers were not considered in our experiments, since they are almost entirely disconnected in the document and hence should be of little interest for most document reconstruction approaches.

In our second experiment we evaluate the number of hypotheses at four stages of algorithm 1, depending on the uncertainty of our simulated text detector. We observe from figure 4 that the number of hypotheses decreases over stages 1–4 regardless of $\sigma$. Stage 1 only counts hypotheses that conform to the estimated orientation of the text detector. The second stage refers to hypotheses retained after the local geometric verification step, followed by stage 3 which retains only the best hypothesis per consensus set. We note that evaluating hypotheses in the initial stage only depends on very few points (from the consensus sets), which makes this local verification very efficient. Finally, we plot the number of hypotheses for stage 4, in which the global error function $\epsilon_2$ defined in eq. (4) needs to be evaluated. As can be seen, only very few (in the order of $10^2$) hypotheses are retained for this last round of selection, after which only the best alignment is retained.

## 7. CONCLUSION

In this work we presented an approach for aligning pairs of hand-shredded document pieces based on boundary information. We discussed how content-based prior information can be used to speed up the alignment process while also improving the system's performance. Furthermore, we showed that constructing consensus sets from boundary points is effective for identifying correct hypotheses. We conclude that our approach needs about four orders of magnitude less global verifications compared to a naive random sampling strategy.

## 8. REFERENCES

[1] J. Schneider and B. Nickolay, "Automatische virtuelle rekonstruktion vernichteter dokumente," *Fraunhofer FUTUR*, vol. 2, pp. 6–8, 2006.

[2] T. F. Smith and M. S. Waterman, "Identification of common molecular subsequences," *Journal of Molecular Biology*, vol. 147, no. 1, pp. 195–197, 1981.

[3] H. Bunke and U. Bühler, "Applications of approximate string matching to 2d shape recognition," *Pattern Recognition*, vol. 26, no. 12, pp. 1797–1812, 1993.

[4] L. Chen, R. Feris, and M. Turk, "Efficient partial shape matching using smith-waterman algorithm," in *CVPR Workshop*, 2008.

[5] A. Stieber, J. Schneider, B. Nickolay, and J. Krüger, "A contour matching algorithm to reconstruct ruptured documents," in *DAGM-Symposium*, 2010, pp. 121–130.

[6] M. Donoser, H. Riemenschneider, and H. Bischof, "Efficient partial shape matching of outer contours," in *ACCV*, 2009, vol. 5994, pp. 281–292.

[7] F. Richter, C. X. Ries, N. Cebron, and R. Lienhart, "Learning to reassemble shredded documents," *Transactions on Multimedia*, vol. 15, pp. 582–593, 2013.

[8] S. Cao, H. Liu, and S. Yan, "Automated assembly of shredded pieces from multiple photos," in *ICME*, 2010, pp. 358–363.

[9] S. Suzuki and K. Abe, "Topological structural analysis of digital binary images by border following," *Computer Vision, Graphics, and Image Processing*, vol. 30, no. 1, pp. 32–46, 1985.

[10] D. Douglas and T. Peucker, "Algorithms for the reduction of the number of points required to represent a digitized line or its caricature," *The Canadian Cartographer*, vol. 10, no. 2, pp. 112–122, 1973.

[11] P. H. S. Torr and A. Zisserman, "Mlesac: A new robust estimator with application to estimating image geometry," *Computer Vision and Image Understanding*, vol. 78, pp. 2000, 2000.

[12] M. Fischler and R. Bolles, "Random sample consensus: a paradigm for model fitting with application to image analysis and automated cartography," *Communications of the ACM*, vol. 24, no. 6, pp. 381–395, 1981.

[13] X. Li, W. Wang, S. Jiang, Q. Huang, and W. Gao, "Fast and effective text detection," in *IPCV*, 2008, pp. 969–972.