# Assembling Paper Fragments using Sparse Presentation and Simulated Annealing Algorithm: A Novel Approach

Ting Wang[1] and Jiansheng Wang[1]

*School of Electronics & Information Engineering,*
*LanZhou Jiaotong University*
*wangting@mail.lzjtu.cn, Wangjsh@mail.lzjtu.cn*

### Abstract

*In this paper, we conduct a novel research on assembling paper fragments with kernel sparse representation and regular edge geometry analysis using the traditional example of rectangular pieces. At the initial stage, we adopt the methodology of image sparse presentation technique to overcome the influence of noise. During the process of assembling, we make good use of MATLAB and C++ to extract the core visual information from the fragments' digitally to capture the matrix in the grey value scale. Edge characteristics are derived and regarded as the basic unit to find out fragments which belong in the first column. According to the similarity characteristic, adjacent rows are found and matched accordingly, annealing algorithm is used to gather the fragments. From the perspective of practical use, we find out the robustness and effectiveness of our proposed approach. Compare with some state-of-the-art algorithms, our methodology shows the better accuracy, it's of great importance to the community of fragment assembly.*

*Keywords: Sparse Presentation, Character Features, Simulated Annealing Algorithm (SAA)*

## 1. Introduction

Automatic fragment assembly as one of the most state-of-the-art research areas on image recovery is frequently studied in the computer science and engineering community. The technique has been broadly adopted including judicial evidence recovery, historical documentation and malfunction analysis. More recently, through restoration of the Stasi files, the more in-depth researches on automatic fragment assembly algorithm have perked wide public interest [1]. The existing research has gained satisfactory advances, in the literature [2], Leitao proposed a novel dynamic programming sequence matching based approach to enhance the robustness. Luo [3] recommended us to use an edge detection algorithm based on link scan line segment fragment of a closed curve edge; Line and circle interpolation method to ensure the continuity of the scanning direction. He [4] suggested an ant colony optimization based algorithm with the combination of contour similarities. Zhao [5] highlighted deficiencies of classic assembly methods using geometry features, by analyzing the text feature, the semi-auto assemble approach is proposed. The above research all target fragments of irregular geometries. With the emergence of paper shredders, fragments have steadily become more regular [6]. In contrast, irregular fragments matching contour features can be assembled and fragment assembly more difficult due to the smooth edges on a regular basis. At present, there are just few researches on assembling regular paper fragments. Zhao [7] extracted pixel numbers from vertical cutting, fragment access to the corresponding matrix on a regular basis. All the existing methods are sensitive to noise, adding de-noising procedure to the present algorithms is indeed needed. Match rate is then set as objective function to assemble fragments by method of exhaustion.

This paper re-assembles rectangular paper fragments cut vertically and horizontally, by digitally extracting the fragments' visual information with MATLAB and C++ software to achieve the corresponding grey value matrix. As an innovative approach, we add the kernel sparse presentation method into the existing algorithm to remove the influence of noise. In addition, on the purpose of finding fragments which belong to the initial column we derive and use the contour characteristics. Adjacent rows are matched subsequently according to similarity and the annealing algorithm (AA) is applied to join the paper fragments by setting edge goodness of fit as the objective function.

## 2. Digitalizing the Fragments

As the paper fragment contours are regular, contour features cannot be used in the assembling process. Digital information of the fragments visual features are extracted and analyzed to achieve reassembly.

Pixel acts as the smallest unit of an image, each image has $m \times n$ cells, $m$ representing length and $n$ width, respectively. The corresponding grey value matrix $A$ has $m$ rows and $n$ columns, $(m, n)$ representing the pixel's position within the image, equivalent to $(x, y)$ in a co-ordinate system. $(m, n) = x_{m,n}$, represents the grey value of $(m, n)$ within the range of 0~255. Each fragment is scanned as an image and we digitalize the image as a $m \times n$ grey value matrix shown in Figure 1. Therefore, we can visually transfer each pixel which is the basic unit of an image into a numerical matrix.
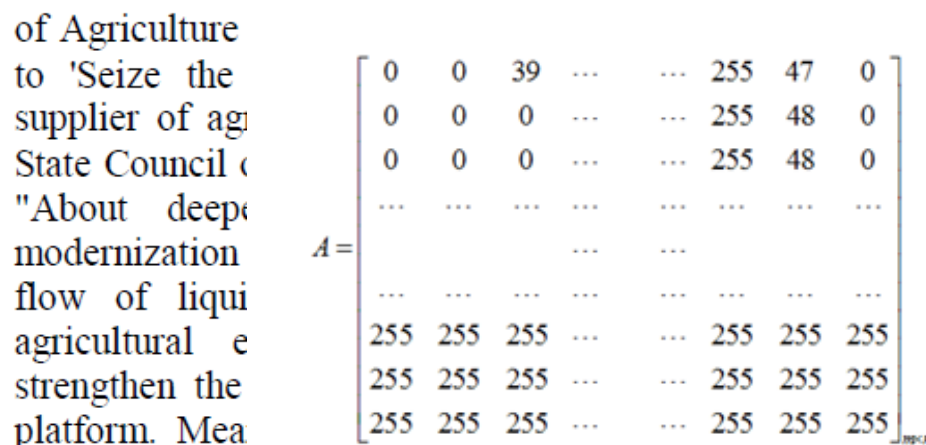


**Figure 1. Grey Value Matrix and Original Image**

## 3. Image Sparse Presentation

### 3.1. Preliminaries of Sparse Presentation

The target of signal sparse presentation(SP) is to use the weighted linear combination of the basic or atom signal [8-11] to better represent the origibal input vecctor $y \in \square^{d}$. We denote the set of basic signal or original data set to be the dictionary $X \in \square^{d \times N}$. As far as the prior knowlege is concerned, the combination coefficient $\beta \in \square^{N}$ is sparse, through mathematical analysis, we could get the sparse coefficient via formula 1.

$$\min_{\beta} \|\beta\|_{0} \qquad st. \qquad X\beta = y \qquad (1)$$

Where, $\left\| \Box \right\|_0$ represent the norm of $l_0$. NP-hard influence the popularity and efficient of the formula (1), therefore, we relat the condition to the following formula.

$$\min_{\beta} \left\| \beta \right\|_1 \qquad st. \qquad X\beta = y \qquad\qquad (2)$$

Which is a Linear Programming (LP) approach and could be solved with less time consum. The literature [12] provides the detailed discussion on the solution to the formula and some applicable implementations. If the problem is just related to the area of remeving noise, we can adopt the Order Cone Programming problem as an alternative.

$$\min_{\beta} \left\| \beta \right\|_1 \qquad st. \qquad \left\| X\beta - y \right\|_2 \le \varepsilon \qquad\qquad (3)$$

Where $\varepsilon$ is a pre-specified error tolerance. In the research area of sparse based image classification, the identity label of the class is denoted as:

$$l(\text{y}) = \arg\min_{j \in \{1,\cdots,C\}} [\text{r}_j(y)] \qquad\qquad (4)$$

Kernel sparse presentation as a state-of-the-art machine learning research area, up-cast the original sparse presentation problem into the space of high-demensional. Suppose there exists a feature mapping function $\Phi : \chi \to F$ which maps the feature and the basis as:

$$\begin{aligned} y &\to \Phi(y) \\ X &\to U = \left[ \Phi(x_1), \Phi(x_2), \cdots, \Phi(x_N) \right] \end{aligned} \qquad\qquad (5)$$

## 3.2. Orthogonal Matching Pursuit (OMP)

Before the compressed sensing(CS) theory was initially proposed, plenty of approaches had been applied for sparse approximation such as wavelet decomposition, timefrequency dictionaries and adaptive timefrequency decompositions. Orthogonal Matching Pursuit (OMP) is one of the effective and applicable approaches. Generally speaking, the orthogonal matching pursuit aims at solving the following problem:

$$\min_{\beta} \left\| \beta \right\|_0 \qquad st. \qquad \left\| X\beta - y \right\|_2 \le \varepsilon \qquad\qquad (6)$$

Given a dictionary $X \in \Box^{d \times N}$, the computational complexity of linear programming is in the scale of $O\left( d^2 N^{\frac{3}{2}} \right)$. As a comparison the computational complexity of OMP is $O(dN)$. In Figure 2, we show the basic steps of OMP.
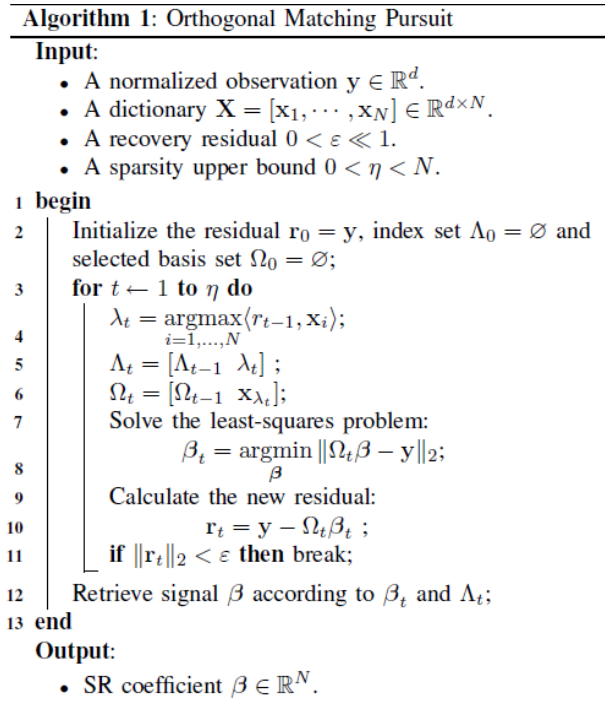
**Algorithm 1:** Orthogonal Matching Pursuit

**Input:**
- A normalized observation $\mathbf{y} \in \mathbb{R}^d$.
- A dictionary $\mathbf{X} = [\mathbf{x}_1, \cdots, \mathbf{x}_N] \in \mathbb{R}^{d \times N}$.
- A recovery residual $0 < \varepsilon \ll 1$.
- A sparsity upper bound $0 < \eta < N$.

1 **begin**
2     Initialize the residual $\mathbf{r}_0 = \mathbf{y}$, index set $\Lambda_0 = \varnothing$ and selected basis set $\Omega_0 = \varnothing$;
3     **for** $t \leftarrow 1$ **to** $\eta$ **do**
4         $\lambda_t = \underset{i=1,\dots,N}{\arg\max} \langle \mathbf{r}_{t-1}, \mathbf{x}_i \rangle$;
5         $\Lambda_t = [\Lambda_{t-1} \;\; \lambda_t]$ ;
6         $\Omega_t = [\Omega_{t-1} \;\; \mathbf{x}_{\lambda_t}]$;
7         Solve the least-squares problem:
8         $\beta_t = \underset{\beta}{\arg\min} \|\Omega_t \beta - \mathbf{y}\|_2$;
9         Calculate the new residual:
10         $\mathbf{r}_t = \mathbf{y} - \Omega_t \beta_t$ ;
11         **if** $\|\mathbf{r}_t\|_2 < \varepsilon$ **then** break;
12     Retrieve signal $\beta$ according to $\beta_t$ and $\Lambda_t$;
13 **end**

**Output:**
- SR coefficient $\beta \in \mathbb{R}^N$.

**Figure 2. The Basic Steps of OMP**

## 4. Assembling Algorithm

Fragment assembly and repair of edge pixels is the standard goodness-of-fit between adjacent segments, better adapt to mean higher probability of adjacency. Calculation process is based on the TSP [13] which transfers the original problem to the combinational optimization issue. Setting the edge pixels of goodness of fit as objective function, and through simulated annealing algorithm, we can automatic assembly pieces of paper. In the flowchart marked Figure 3, we present the approach.

### 4.1. An Overview of the Algorithm

Simulated annealing algorithm [14-19] is a random optimization algorithm based on the Monte Carlo method; it is a combination of the metallurgical annealing process and combinational optimization. Through random search in the preset parameters, the algorithm could avoid local optimal probability under the global optimal solutions and methods. Process of the Simulated annealing algorithm is shown in Figure 4.
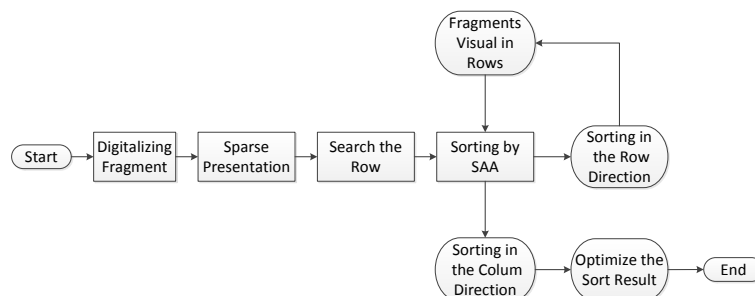
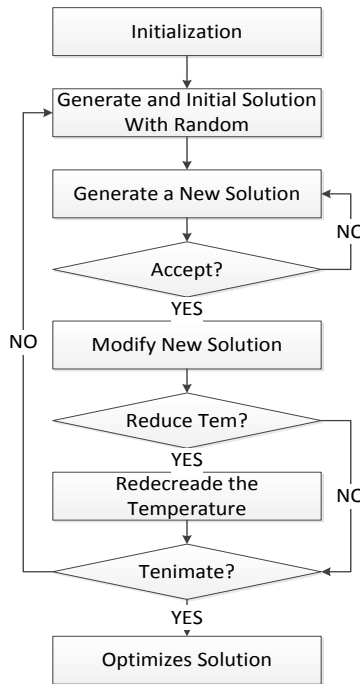**Figure 3. Auto Assembly Process Overview**

**Figure 4. Process of Simulated Annealing Algorithm**

## 4.2. Determining Edge Fragments of the Original Document by Row

In this study, the original document is separated into 11 rows and 19 columns, thus rectangular fragments as shown in Figure 5. By searching the white edge, the edge of the original document fragment can be the first place. 255 examples, looking for gray left edge value, each pixel of image matrix, the matrix elements of the search range is found in the left side, if the sum of matrix is 255 multiples, moments have a blank left edge. Through this method, it was found that there are exactly 11 fragments each for the far left column 1 and far right column 19. It was discovered that there were 22 and 35 fragments which satisfied the criteria for having a blank top edge and blank bottom edge respectively, contradicting the reality of 11 fragments on each edge. This means the above search method does not yet suffice, and further searches are required according to text line spacing. 1 fragment from 11 of those in column 1 is analyzed with line spacing and its relative vertical position within the fragment as bounding condition, filtering fragments in the same row. As shown in Figure 6, three fragments have three identical line spacing, and identical positions within the fragment. By searching for similar information, the 19 fragments of that row can be located. The same method can be applied to locate the remaining 10 rows.



**Figure 5. A Sample Document Page**

of Agricult    ure and th    e Ministry oi
to 'Seize       the comm    anding heig
supplier of ? agricultu    al products';
State Coun    cil of the (   CPC Central

**Figure 6. Three Sample Fragments in the Same Row**

### 4.3. Search Process for Fragment Rows

According to the above analysis, identify the serial number of the segments of each method in response to the following procedure: Step 1: Selecting a fragment randomly within the first 11, and find fragments of the same row by using its line spacing information. Step 2: Find the sum of pixel point values of each row within an image to get a column matrix. Step 3: Determine if the 180 numbers of the column matrix can be divided completely by 255, if so record as 1, if not, as 0, thus establishing a 0-1 column vector sized . Step 4: Find sum of the remaining 208 pixel point values and determine whether each value can be divided completely by 255, establishing a 0-1 column vector. Step 5: From the far left 0-1 column vector, randomly select 2 ~ 5 sections of consecutive 1 values, as the defining region for the row's line spacing. Other fragments of the same row are found through defining those images with the same defining regions in the same positions. Step 6: Randomly select another fragment from the first column, repeat steps 2~5 until 19 fragment serial numbers are found for each row.

### 4.4. Fragment Assembly by Row

Four corners pieces fragments sequence can be found by searching the advantage. Looking for the upper left corner, the top and left edge is blank, the serial number should be crossing the left-most column and row element. Also, the other three Angle can be found. Serial number of the four corner pieces exist within each column and row sets, so it can identification from the edges of the clip series. Arrange 19, the pieces of paper each line using simulated annealing algorithm. For two adjacent fragments calculate the logical value of pixel matrix $S_{i,j}(k,72)==S_{i,j+1}(k,1)$, ($k=1,2,\cdots,180; i=1,2,\cdots,11; j=1,2,\cdots,18.$), regarded to be the rate of matching for these fragments. As far as all fragments are assembled, we will hereby gain the better accuracy if the matching rate is enhanced. Figure 7 and Figure 8 are partial edge information of adjacent fragments within a row. The serial number of the implementation of the method to determine each fragment in the continuous by simulated annealing algorithm are as follows:

Step1: From the above search results, extract serial numbers of a fragment row. The first and last fragments are already identified.

Step2: Use the Monte Carlo iteration to generate an initial feasible solution index $x(0)$ as the initial solution.

Step3: Initiate the programme from the current temperature $T_i$.

Step4: Set match rate of the right most edge pixel value on fragment $j$ and left most edge pixel value on fragment $j+1$ as $F$, and set the objective function to maximize the global match rate. Calculate the objective function value $F_0$.

$$M \text{ ax } F = \sum_{j=1}^{18}\left(S_{i,j}(k,72)==S_{i,j+1}(k,1)\right) \quad k=1,2,\cdots,180; i=1,2,\cdots,11; j=1,2,\cdots,18 \tag{7}$$

Step5: Use the initial solution as core to create stochastic disturbance within the solution space. Generate a perturbation solution $x'$, and calculate its objective function $F'$.

Step6: Determine if the obtained perturbation solution satisfies formula (8). If the generated function value $F'$ of $x'$ is greater than $F_0$ of $x(0)$, then accept $x(1) = x'$ as a new solution of this iteration; otherwise, accept $x'$ as the new solution with a probability of $e^{\frac{f(x')-f(x(0))}{T}}$ .

$$G(F_0 \rightarrow F') = \begin{cases} 1 & f(F') < f(F_0) \\ e^{-\frac{f(F')-f(F_0)}{T_0}} & \text{Other} \end{cases} \qquad (8)$$

Of which, $x(k), k = 0,1,2...n$ is the initial solution, $x'$ the perturbation solution.

Step7: After running multiple iterations in temperature, cool by mechanism

$$T' = \alpha \times T \qquad (9)$$

Of which, $\alpha \in (0,1)$, stipulating $\alpha = 0.999$, $T$ is the current temperature, $T'$ is the temperature of the next cooling phase, $T_0 = 1$ is the initial temperature.

Step8: Determine if the termination temperature $e = 10^{-20}$ is satisfied, if $T \le e$, terminate the cooling mechanism and algorithm. If not, repeat steps 3~7. By optimizing the arrangement of 19 fragments in each row, the near-optimum solution for each row can be found. Figure 9 is the assembled image of a fragment row, thus is the implemented process of assembling paper fragments by row using our proposed methodology.



**Figure 7. Two Fragments in the Same Row**



**Figure 8. Marginal Grey Value Matrix of Adjacent Fragments**



**Figure 9. Auto Assembing Result of a Row**

### 4.5. Fragment Assembly by Column

The assembly of paper fragments by column is similar to that by row discussed priously, the process also implemented through simulated annealing algorithm. When assembling the columns have more edge pixels, this process is more accurate and more efficient. After completion of row and column assembly, we can restore the original document. Though the assembly process in this study used document with line spacing as example, the process can also be used for documents without line spacing.

## 5. Conclusion and Summary

In this research paper, we propose a novel methodology on assembling paper fragments with regular edge geometry and sparse presentation. Through scanning the paper fragments we have it digitalized and capture the visual information with the analysis of grey value matrix and corresponding features. Line spacing and spacing position is set as constraints to find matching fragments within the same row; the edge pixel match rate are set as objective function to find the optimum solution through simulated annealing algorithm. Fragments are assembled within each row, and then vertically row by row.

The study not only considering the edge pixel matching rate is visual characteristics amplified assembly precision. Simulated annealing procedure is used to convert the challenge of automatic assembly into a set of optimization, improve the efficiency of assembly through macro management. However, it is important to note that in the case of crushing, proposed method are not able to guarantee high precision because of the lack of visual information in each fragment, therefore, in turn, will affect the assembly efficiency. In the future research, we plan to use some mathematical optimization method to modify our proposed algorithms so that the robustness and accuracy of the methodology will be enhanced rapidly.

## Acknowledgment

## References

[1] L. da Gama, H. Cristina and J. Stolfi, "A multiscale method for the reassembly of two-dimensional fragmented objects", Pattern Analysis and Machine Intelligence, IEEE Transactions on, vol. 24.9, **(2002)**, pp. 1239-1251.

[2] H. C. G. Leitao, R. F. V. Saracchini and J. Stolfi, "Matching photometric observation vectors with shadows and variable albedo", Computer Graphics and Image Processing, **2008**. SIBGRAPI'08. XXI Brazilian Symposium on. IEEE, **(2008)**.

[3] Q.-X. Huang, "Reassembling fractured objects by geometric matching", ACM Transactions on Graphics (TOG), ACM, vol. 25, no. 3, **(2006)**.

[4] F. Amigoni, S. Gazzani and S. Podico, "A method for reassembling fragments in image reconstruction", Image Processing, ICIP Proceedings International Conference on IEEE, vol. 3, **(2003)**.

[5] J. C. McBride and B. B. Kimia, "Archaeological fragment reconstruction using curve-matching", Computer Vision and Pattern Recognition Workshop, **2003**. CVPRW'03. Conference on IEEE, vol. 1, **(2003)**.

[6] G. Papaioannou, E-A. Karabassi and T. Theoharis, "Reconstruction of three-dimensional objects through matching of their parts", Pattern Analysis and Machine Intelligence, IEEE Transactions on, vol. 24.1, **(2002)**, pp. 114-124.

[7] N. Memon and A. Pal, "Automated reassembly of file fragmented images using greedy algorithms", Image Processing, IEEE Transactions on, vol. 15.2, **(2006)**, pp. 385-393.

[8] H. Wang and J. Wang, "An Effective Image Representation Method using Kernel Classification".

[9] A. C. James, *et al.,* "Spatially sparse pattern-pulse stimulation enhances multifocal visual evoked potential analysis", Investigative Ophtalmology and Visual Science, vol. 46.5, **(2005)**, pp. 3602.

[10] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine", The journal of machine learning research, vol. 1, **(2001)**, pp. 211-244.

[11] H. Zou and R. Li, "One-step sparse estimates in nonconcave penalized likelihood models", Annals of statistics, vol. 36.4, **(2008)**, pp. 1509.

[12] B. Cheng, "New existence and multiplicity of nontrivial solutions for nonlocal elliptic Kirchhoff type problems", Journal of Mathematical Analysis and Applications, vol. 394.2, **(2012)**, pp. 488-495.

[13] H. Zhu and Y. Zhong, "Computer technology and development", vol. 19, no. 6, **(2009)**, pp. 32-35, in Chinese.

[14] G. Reinelt, ORSA journal on computing, vol. 3, no. 4, **(1991)**, pp. 376-384.

[15] K. L. E. I. N. Bouleimen and H. O. U. S. N. I. Lecocq, "A new efficient simulated annealing algorithm for the resource-constrained project scheduling problem and its multiple mode version", European Journal of Operational Research, vol. 149.2, **(2003)**, pp. 268-281.

[16] Y.-J. Jeon, *et al.,* "An efficient simulated annealing algorithm for network reconfiguration in large-scale distribution systems", Power Delivery, IEEE Transactions on, vol. 17.4, **(2002)**, pp. 1070-1078.

[17] I. Dupanloup, S. Schneider and L. Excoffier, "A simulated annealing approach to define the genetic structure of populations", Molecular Ecology, vol. 11.12, **(2002)**, pp. 2571-2581.

[18] T. K. Varadharajan and C. Rajendran, "A multi-objective simulated-annealing algorithm for scheduling in flowshops to minimize the makespan and total flowtime of jobs", European Journal of Operational Research, vol. 167.3, **(2005)**, pp. 772-795.

[19] T.-H. Wu, C.-C. Chang and S.-H. Chung, "A simulated annealing algorithm for manufacturing cell formation problems", Expert Systems with Applications, vol. 34.3, **(2008)**, pp. 1609-1617.

## Author

**Ting Wang**, she received his M.Sc. in Lanzhou Jiaotong University (2008). Now she is full instructor of communication at Internet of Thing Department, University. Since 2013 she is Member of CCF. Her current research interests include different aspects of Artificial Intelligence and Distributed Systems.