



HAL
open science

User-driven Page Layout Analysis of historical printed Books

Jean-Yves Ramel, Marie-Luce Demonet, Sébastien Busson

► **To cite this version:**

Jean-Yves Ramel, Marie-Luce Demonet, Sébastien Busson. User-driven Page Layout Analysis of historical printed Books. *International Journal on Document Analysis and Recognition*, 2007, 9 (2-4), pp.243-261. 10.1007/s10032-007-0040-6 . hal-00150167

HAL Id: hal-00150167

<https://hal.science/hal-00150167>

Submitted on 2 Jun 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

User-driven Page Layout Analysis of Historical Printed Books

J.Y. Ramel*, **S. Leriche***, **M.L. Demonet****, **S. Busson****

* *Lab. d'Informatique, Ecole Polytechnique de l'Université de Tours,
64, avenue Jean Portalis 37200 Tours- France
Tel : +33.2.47.36.14.26 Fax : +33.2.47.36.14.22
ramel@univ-tours.fr*

** *CESR - UMR 6576 du CNRS –
59, rue Néricault-Destouches - BP 11328 - 37013 Tours-France*

ABSTRACT. In this paper, based on the study of the specificity of historical printed books, we first explain the main error sources in classical methods used for page layout analysis. We show that each method (bottom-up and top-down) provides different types of useful information that should not be ignored if we want to obtain both a generic method and good segmentation results. Next, we propose to use a hybrid segmentation algorithm that builds two maps: a shape map that focuses on connected components and a background map which provides information about white areas corresponding to block separations in the page. Using this first segmentation, a classification of the extracted blocks can be achieved according to scenarios produced by the user. These scenarios are defined very simply during an interactive stage. The user is able to make processing sequences adapted to the different kinds of images he is likely to meet and according to the user needs. The proposed “user driven approach” is capable of doing segmentation and labeling of the required user high level concepts efficiently and has achieved above 93% accurate results over different data set tested. User feedbacks and experimental results demonstrate the effectiveness and usability of our framework mainly because the extraction rules can be defined without difficulty and parameters are not sensitive to page layout variation.

KEYWORDS: page layout analysis, segmentation, digital libraries, analysis strategy, zone classification

User-driven Page Layout Analysis of Historical Printed Books

1. Introduction

High level analyses of document images are mainly based on the output of a page segmentation process. For example, the extracted text regions can be the input to an OCR system to retrieve the ASCII characters printed on the pages. The spatial relationships between segmented blocks along with other features can be used in logical page organization analysis to group the extracted components appropriately and recover the correct reading order. Many techniques for page segmentation have been proposed in the literature [Baird90][Le01][Hadjar02][He03][Marinai05] but most of them are based on the assumption that an input document image consists of a set of rectangular blocks. Furthermore, the classification step is generally domain specific and uses static rules to automatically determine, for each block, the coherent label selected from a predefined list (title, paragraph, graphic, table,...). These limitations appear too restrictive with respect to historical documents and new approaches need to be developed.

In this context, we present a work achieved in collaboration with the “Centre d’Etude Supérieur de la Renaissance” of Tours (CESR / <http://www.cesr.univ-tours.fr>). The CESR is a training and research centre which receives students and researchers who wish to work on various domains of the Renaissance using a rich library of historical books. The CESR wants to create a Humanistic Virtual Library; however, until now, only bitmap versions of several books that have been scanned or photographed are accessible. The initial objective of the CESR was to obtain an ASCII version of the text contained in the pages of these historical books. The centre first tried to use the commercial OCR software to index their books but they quickly realized that, applied to historical documents, this procedure would have been vowed to failure. So, the CESR asked our Pattern Recognition and Image Analysis research team to help them to define a new system adapted to their needs. They have appreciated our efforts as our collaboration will lead to a system

able to bring a better description and indexation of the content of their books and would also make the search and the reading of these precious historical books easier.

In this paper, we first study the old book specificities in order to infer some invariant characteristics used during the automatic analysis of documents with complex structures. Next, we describe the layout extraction methods applicable to such documents by focussing on their quality and their drawbacks. We also specify the adaptations we need to make in the classical methods in order to analyse complex documents. The second part of the article proposes a new hybrid method for the extraction of the layout of the documents based on the construction of two representations of the contents of the images. A mapping of the shapes and a mapping of the background are computed. By exploiting this information, our algorithm produces and sends back a list of blocks constituting a first segmentation result.

Then, this initial representation of the image is used during a more sophisticated analysis. Having an aim of genericity, the architecture of the system that we carried out authorizes an interactive installation of scenarios for analysis of the image contents. Scenarios work on the initial representation provided by the first step of the segmentation. According to its needs (localization of the ornamental letters, the notes at margins, titles,...) and using user-friendly interfaces, the user (not expert in image processing) builds scenarios allowing to label, to merge, to remove the blocks contained in the intermediate representation. One can thus locate the desired entities without taking care of the other areas of the image. The elaborated scenarios can then be stored, modified and applied to various sets of images during batch processing.

2. A brief study of the characteristics of historical books and of processing methods

2.1. *Characteristic of historical books*

2.1.1. *Production values*

Since its creation, the CESR has established a collection of precious historical books currently numbering round about 3000 copies, dating from the middle of the XIV century to the beginning of the XVII century. The first books belong to the beginning of the printing era when the fonts used and the layouts of the pages were very close to those of the handwritten books. The collection of the CESR is pan European: coming from France, Germany, Italy, Switzerland, and Holland. The languages used, Latin or French, bring an additional factor of variability to the books. The Renaissance typographies used have significant variability. Some examples of images of these precious books are presented in Figure 1.

The technical materials used impose particular presentations in the layout of the pages. The variability of page layouts is much larger than in the current books due either to inaccuracies or to liberties taken by the printer. Most of the time, the text part represents the majority of the page, with notes in the margins on either side of the text. The page can also contain graphical parts of various sizes and some *ornamental patterns*. In the text, we can find known structures like the titles and the subtitles, the paragraphs, the page numbers, and other more particular structures like the *catchwords*. These are the words which are placed at the bottom of the page, they represent page number rather they act as a connection between two pages and usually corresponds to the first word of the next page. The catchwords help the printers to organize the manuscript correctly. The styles used can alternate, normal style, justified or aligned on the left. Another characteristic of old printed books comes from weak separations between blocks of text (notes at the margins and body text for example). Lastly, on some pages current page layout rules are not complied with: for example an illustration can overflow into the margins (figure 1). In the Renaissance books, the illustrations were generally printed using metal or wood plates, engraved and inked to reproduce the image. They are generally included in a white rectangle often surrounded by text.

2.1.2. Characteristics due to image acquisition techniques

Other characteristics are due to the digitalization processes. They relate to the lack of lighting in the binding, to the curve of the text lines, to page distortion, and to the imperfect elimination of the stains. Much research was undertaken, in particular at the time of the Debora project, to correct these defects [Trinh03] [Lebourgeois03].



Figure 1: Examples of pages of historical books illustrating close spaces between blocks, complex

layouts with texts and graphics at random positions

2.1.3. Selected invariant characteristics

This study enables us to extract a list of characteristics which are essential to take into account while designing algorithms for layout analysis of the complex documents. Here is the prior knowledge that we used to implement our algorithms:

o Characteristics depending on printing materials and layout conventions of the Renaissance:

- Complex page layout which can present several columns with irregular sizes
- No Editorial Style or identifiable logical structure

- Presence of printed or handwritten notes at the margins
 - Presence of location indicators: line numbers, page numbers, catchwords ...
 - Use of specific and multiple fonts
 - Frequent use of ornaments (non textual blocks) such as borders, ornamental letters...
 - Variable location of the graphical illustrations and the associated legends
 - Absence of leading bringing contacts between characters
 - Non constant spaces between characters, words and blocks of text
 - Text blocks are not rectangular
- Image characteristics connected to the nature of the documents:
- Images still degraded even after restoration (appearance of the characters on the back of the page)
 - Presence of superposition of information layers (noise, handwritten notes...)

2.2. Evaluation of existing methods on images of historical documents

Typography and printing technology have progressed enormously and current books meet quite a different standard of presentation than historical books. Consequently, the software built to recognize current documents is often inadequate when processing books of the Renaissance period. The methods of layout extraction employed by this software can be classified in 3 main categories: bottom-up, top-down and hybrid methods [Belaid97].

- To be classified among the bottom-up methods, a method must work starting from the image material in order to obtain a more and more interpreted structure using the information provided by the image rather than information about the model of the document. So, most often, such data-

driven methods work using an element aggregation strategy. In this category, we find the methods based on different morphological or filtering techniques, methods based on texture analysis and methods based on the study of the connected components in the image. Thus, the problems involved in the bottom-up methods come from many parameters which are difficult to set and from prohibitive processing times. Moreover, these methods do not allow the analysis of noisy images of documents.

- In contrast, the top-down approaches use knowledge about the document to drive the segmentation of the image (model driven methods). The segmentation process is very different from the one used by the previous methods because algorithms look for precise information specified in the model of the documents. So, the image analysis is achieved by using mostly a split strategy forced by a priori knowledge. Most of the attempts for creating such systems have come to very complex architectures or to applications adapted to a restricted set of documents. Nagy et Al [Nagy93] have, for example, developed a syntactic based system working with technical journals. They define a set of appropriate context free grammars, each defining rules to locate more and more structured entities till up to the logical objects. From the grammars, parsers are automatically obtained: they are then used to perform segmentation and labelling simultaneously. A set of alternative grammars can be used to allow different document structures to be extracted and checked.
- The hybrid approaches try to use advantages of both top-down and bottom-up methods. Such systems generally include a set of images analysis and page segmentation algorithms that work without requiring any prior information about a specific model of a document. In contrast, a controller has to be defined to select which algorithm should be run at a given moment according to rules which are valid within a specific category of documents.

Document analysis systems of these two last categories employ, often in an implicit way, a set of typesetting rules and have been designed to operate in specific domain mainly due to the test images used.

This fact gives rise to a variety of methods which are characterized by different objectives and different assumptions about input type, layout structure and expected data types. A weak aspect which is common in all these algorithms is the lack of information about the handling of regions which do not belong to one of the expected categories.

Thus, before describing our system, we present the results of some experiments which enabled us to evaluate the various traditional methods used for segmentation of traditional documents on the historical printed documents.

2.2.1. A few words about pre-processing

Despite the many efforts spent on the development of skew estimation algorithms, every year new algorithms are proposed in the literature. This is mainly due to the need of both accurate and computationally efficient algorithms, and of methods that do not make strong assumptions about the specific class of documents they can deal with. Today, commercial solutions exist and are even regarded as satisfactory tools (Book Restorer [Digibook2005]) although some problems still remain. For example, the correction of the page distortion (for example the skew) can involve degradation on the borders of the text and the images areas. This is probably due to the general assumption that text represents the most relevant part of the document image. So, performance decreases due to the presence of other components like graphics or pictures. Furthermore, the major part of the algorithms assumes dealing with documents with a clearly dominant skew angle, and only a few methods can deal with documents having multiple skew angles [O’Gorman93][Le94] [Spitz98] [Min96]. In this work, we decided not to spend time on working with these low level problems but rather to set up algorithms able to deal with these specificities.

2.2.2. Morphological and differential Filtering algorithms

Morphological methods have been implemented and frequently tested mainly for the fusion of characters in words, then in lines and for the elimination of the noise [Chen95] [OGorman95]. Figure 2 illustrates the results we can obtain using this type of method. The traditional Run Length Smearing

Algorithm (RLSA) [Wong82] can also be used; it works and provides similar results to all those morphological methods.

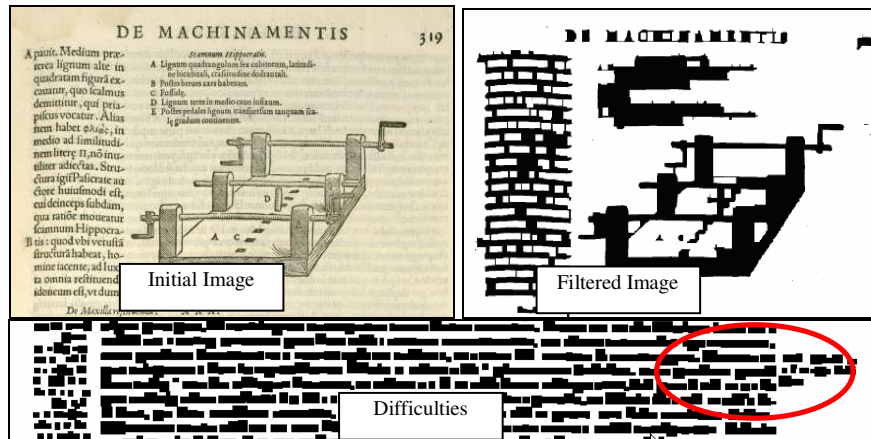


Figure 2: Use of morphological filtering or RLSA is not sufficient on historical books because of horizontal and vertical proximity between different text blocks

It is also possible to work directly on grey level images by using differential filters [Lebourgeois99]. The idea is to use filters that detect and agglomerate the variations of intensity periodically produced by contours of the characters and then to look for horizontal alignments corresponding to text areas.

2.2.3. Methods based on the analysis of the connected components

These methods transform each page into a set of connected components [Wang95]. Figure 3 shows that the size, the proximity, and the relative position of the connected components can be used to extract the physical structure of a document. The bounding boxes of the connected components frequently overlap themselves in the graphical parts but not so much in the text parts. Furthermore, the graphical areas correspond to the connected components whose dimensions (width or height) exceed a fixed threshold.

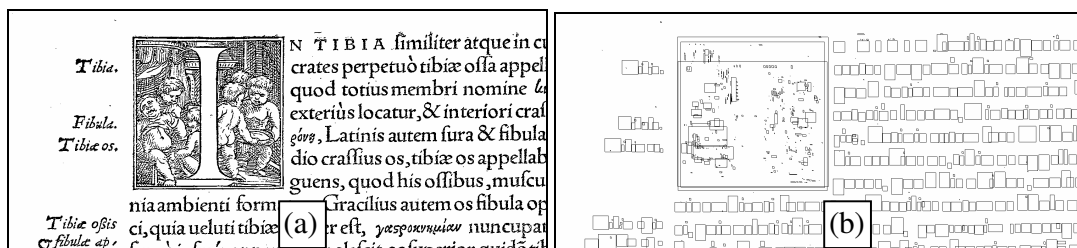


Figure 3: Use of connected components for segmentation

(a) initial image, (b) obtained connected components

A text block is a set of small connected components in close proximity to one another. Two characters are known as neighbours if the distance separating them is lower than a maximum threshold. O' Gorman proposed to use only the neighbourhood between components to locate the Text areas in an image [OGorman93]. Considering the same idea, we implement a method which seeks for the components in the four principal directions of for each text component. A component which does not have a neighbour in at least one direction is considered as a part of the contour of the block. Then, a closed circular chaining translated the border of a text area (figure 4).

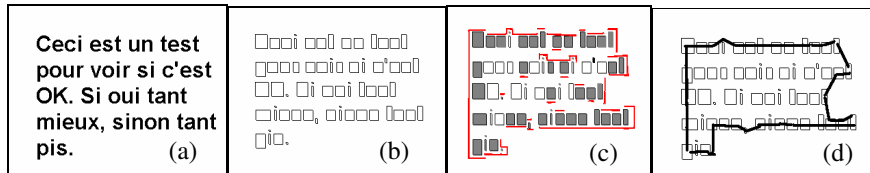


Figure 4: (a) Initial image, (b)connected component, (c) in grey: components of contour, in red: sides without neighbour, (d) circular chaining obtained

K. Kise [Kise98] proposes to use an area Voronoï diagram instead of directly connected components to extract text blocks. The area Voronoi diagram enables to build a neighbour-graph which represents neighbours of the characters. Then, the task of line extraction is considered to be a path selection problem in the graph.

Like the methods based on filtering, these methods are not suitable for historical books processing because of the proximity between some of the text areas in the pages. For the localization of the graphical parts, the thresholds on the size of the connected components must correspond to a size slightly larger than the greatest character contained in the book. However, even if statistical studies can make it possible to automate the selection of these thresholds, these methods hardly support variability in size of the characters and variability in page layout which can frequently appear in the old books. The bad quality of the images (holes, spots...) can also raise many problems when using such methods.

In conclusion, it seems quite difficult to segment correctly the text paragraphs and the graphical parts in the old documents by using bottom-up approaches mainly because they are using only local criteria.

2.2.4. Background analysis

Some other methods are based on the assumptions that the white spaces act as separators of the different regions in the document [Antona98] [Breuel03] [Baird90]. In practice, they consist in a hierarchy of constraints on the relative size of spacings between different objects, e.g. spacings between characters are smaller than spacings between words, which, in turn, are smaller than spacings between lines and so on. In [Breuel03] and [Baird90], the search for the white separations could be achieved using horizontal and vertical projections. Then, the localization of white separations is done by an analysis of the histogram of the number of black pixels of the lines and the columns of the image. Next, one can regard for a significant difference between successive values in the histogram as delimitation between two blocks of text. The problem is to estimate the significant differences between the histogram values. Generally, a priori knowledge about the document to be analysed (number of columns, margins...) is used to locate separations more easily. In the old documents, the tests we have done using such methods, showed that the quality of the results is increased by blackening the bounding boxes of the connected components of the image before calculating the histogram.

The main drawback of these two methods is that they suppose the existence of rectangular text blocks. It is not the case for historical books. Moreover, the problem of this type of techniques relates to the separation between the local and global minima in the histogram. In historical books, certain separations between blocks correspond to non significant minima (case of badly deskewed images). So, it is thus difficult to use these methods when page layout is variable or when the pages are badly rectified by the pre-processing.

The method described in [Antona98] attempts to try to locate in a more local way the inter-block spacings instead of the maximum size white rectangles themselves. It used white tiles to extract contours of regions. But this always supposes that, two blocks are necessarily separated by a white space from

significant surface in either a horizontal or a vertical direction. Unfortunately, as shown in figure 1, it is not the case with the historical documents which we have to process.

Among all the top-down methods suggested to segment complex page layouts of composite documents (newspapers, pages of magazines...), we decided to adapt the Split and Merge method suggested by Hadjar [Hadjar01] during the "segmentation contest" of ICDAR' 01. Its goal is to locate the homogeneous white blocks by splitting the image in smaller and smaller areas (figure 5). A priori information is needed particularly to define the stop criterion in the recursion. This method can be compared with other more or less advanced techniques like the method *Recursive XY-Cut (RXYC)* [Nagy84] [Akindele93]. Readers interested in such method can also look at [Wang89] where Wang and Srihari compare the RLSA and RXYC approaches. In all the cases, it is the analysis of the neighbourhood (realized with the assistance of graphs or trees) of each white area discovered by splitting that will make possible to locate and characterize the blocks contained in the page in a more or less fine manner. The major problem of such techniques is that they are suitable only for layouts that are decomposable by a sequence of horizontal and vertical subdivisions.



Figure 5: Use of the background for page segmentation of historical books
(with a modified version of Split & Merge algorithm)

The top-down methods seem to be less sensitive to the noise than the bottom-up methods. When the images are well deskewed, they partly solve the problem of proximity between blocks and characters. Nevertheless, they are not easily applicable on the documents having a variable or a non rudimentary page layout since they require the use of a priori knowledge [Hadjar02] about the document (number of columns, splitting criteria ...). They are thus not adapted to the processing of the old documents which do not respect a precise Editorial Style.

2.2.5. Multi-resolution and texture analysis

Jain and Al [Jain92] [Jain96] present a multi-channel filtering approach to texture segmentation. The basic assumption is that regions of text in a document image define a unique texture which can be easily captured by a small number of Gabor filters. As said in [Lee01], one major problem of texture-based approaches is their high time complexity since different filters are needed to capture the desired local spatial frequency and orientation of the regions. Many masks are needed to extract local features and small masks do not allow to detect large-scales texture. To avoid this problem, some researches [Eternad97] have proposed to use a multi-scale analysis. Such methods improve the performance but are quite time-consuming due to the dependence on thresholds and parameters. Furthermore, in cases of historical documents, different regions can be incorrectly merged always due to proximity between different text blocks.

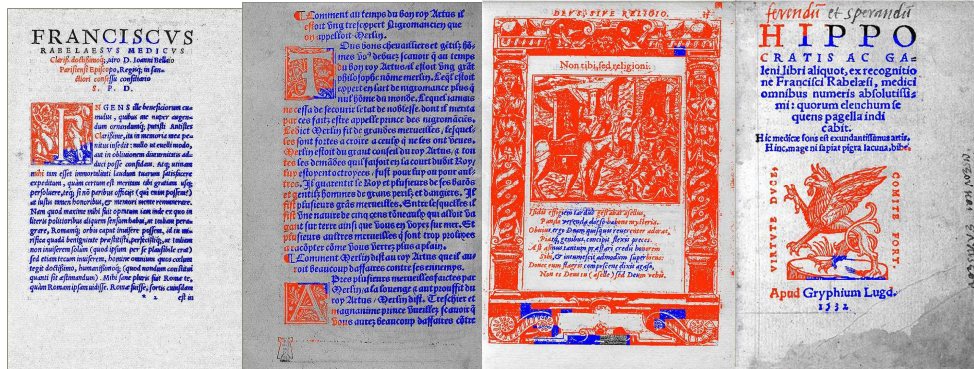


Figure 6: Examples of the results obtained with a texture-based approach (Text regions are in blue and Non-text are in red)

To test the texture-based approaches, we implemented a method based on the extraction and the quantification of the various orientations present in the different areas of the document images based on the use of the autocorrelation function [Journe05]. The conclusion of our experiments (Figure 6) is that Text/Non-text segmentation in Renaissance documents can be viewed as a texture segmentation problem. We also conclude that, the above methods consisting in analyzing the various textures contained in the image, even with a multi-scale approach, are not easily adaptable to obtain the precise contours of the text blocks when the text blocks are too close one to another. Small white spaces become undetectable when the resolution decreases and two different text blocks do not always present any significant differences in their texture. These techniques are nevertheless interesting when we try to characterize globally pre-segmented blocks using particular indices.

3. A new method of segmentation

The preceding tests stressed the limits of the traditional methods to segment historical documents as well as the reasons of their failure. On this basis, we developed a new method exploiting the benefits of the top-down methods and the bottom-up methods at the same time. Our method uses a map of the background of the images to highlight the separation between blocks and a map of the shapes present in the image (foreground). Then, we propose to use simultaneously the information provided by these 2

representations (foreground and background) to segment the image. Thus, we solve most of the difficulties mentioned in the previous chapter.

3.1. Map of the foreground (of the shapes)

The connected components provide relevant information about shapes (graphical and text parts). They correspond to one or more characters, noise, or graphical parts. Their positions, their size, the overlapping of their bounding boxes provide precise information about the contents of the pages. They can provide local information on each shape present in the image. The connected components representing the letters of a word are extremely close. Unfortunately, in old documents, the page layout and the spacing between shapes may change frequently. For example, the last letter of a line can be closer to a note in the margin than to the letter which precedes it on the line. So this information should not be used solely.

To obtain the map of the shapes, after a binarisation of the image (using the method described in [Sauvola00]), we carry out a contour tracking of the shapes which allow us to extract the bounding box from each component. The position and the size of the rectangles are stored in a list which correspond to the foreground map (figure 7a). According to its dimensions, each shape is labelled using one of the following labels:

- *Noise* (connected components of small size)
- *Graphic* (connected components of large size)
- *Text* (connected components of intermediate size)

The thresholds used during this phase are chosen by the user according to the maximum and minimum size of the characters in the book. As discussed in section 2.2.3., the methods based on connected components while using simply their dimensions could not fulfill the whole task. In fact, in our approach, this step provides only a first hand labeling which will be checked and evolved thereafter during the interactive analysis. Our experiment on old documents shows that it is sufficient to separate text shapes

from graphics but to extract and label precisely different parts of a document (paragraphs, illustration...) more elaborated processing steps are needed. These processing steps will be described in section 3.3.

3.2. *Map of the background (separations between blocks)*

In the bilevel image, the background of the page is represented by white pixels. Normally, a great number of white pixels can be aligned vertically or horizontally in the regions between two blocks of a page (text, graphics...). In the same way, a large number of white pixels are aligned horizontally in the regions between two paragraphs. In contrast, the number of white pixels that we are able to align vertically or horizontally in a paragraph between two letters of a word and between two words of a sentence are small, the same between two lines of the same paragraph.

Consequently, we propose to associate to each pixel $P(i,j)$ of the image, the summation of $lgb_h(i,j)$ the number of successive white pixels on the horizontal segment containing the current pixel (line) with $lgb_v(i,j)$ the number of successive white pixels on the vertical segment containing the current pixel (row) to build a distance map. The maximum value for $Ng(i,j)$ is $Max = Height \times Width$ (size of the image) and if necessary the lgb_h and lgb_v values can be weighted respectively by the width and the height of the image in order to not privileged one of the two directions. This summation corresponds to the formula (1) that allows to assign a value to each pixel in the background of the page (white parts in the map of the shapes in which the interior of the connected components are blackened).

$$Ng(i, j) = Max - [lgb_v(i, j) + lgb_h(i, j)] \quad (1)$$

After a normalisation to 255 of $Ng(i,j)$, we obtain as shown in the figure 7b a map with grey levels). As the figure 7b shows it, our map of the background translates the separation (more or less underlined) between the blocks of the page.

Baird [Baird92] and Antonacopoulos [Anto98] have proposed a method based only on information similar to that provided by our map since they proposed to extract the white areas from maximum size in a

page. But, because of the use of rectangles or of tiles, these methods are more sensitive to the distortion and noise problems than the one proposed here. Moreover, in our case, the white blocks constitute only one part of the information which we exploit to carry out the segmentation (and will be combined with the foreground map). In the same way, compared with methods based only on the analysis of the neighbourhood of connected components like [Kise98], our approach uses a more global information because each value in the background map depends on the global layout of the page.

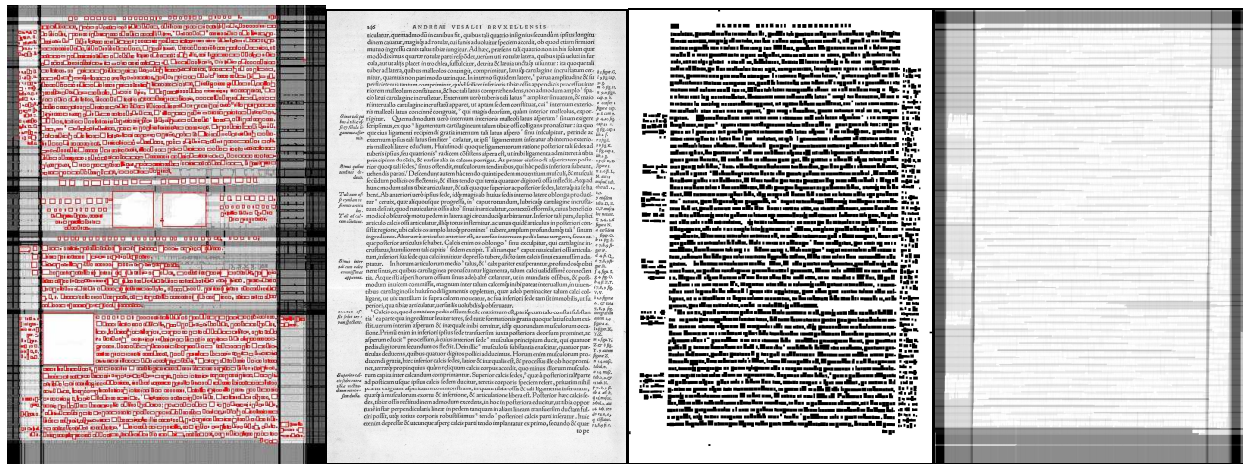
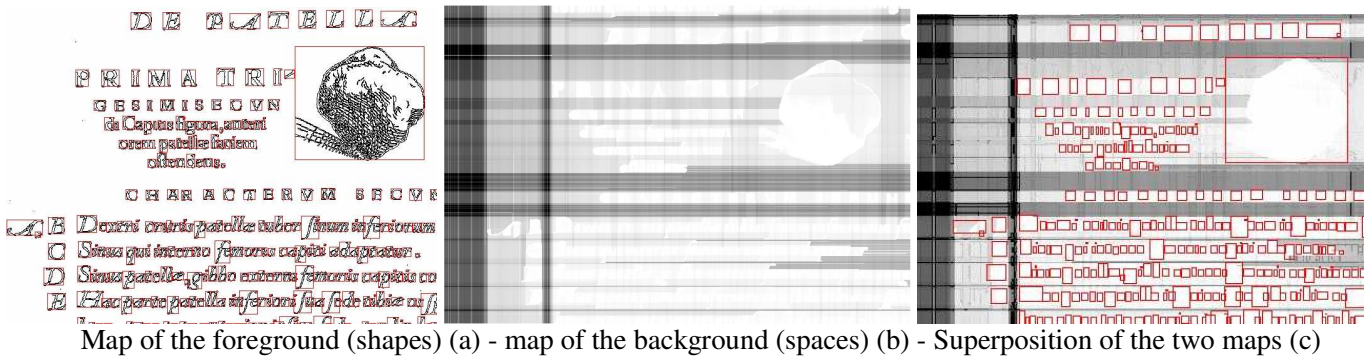


Figure 7: Maps of the background and shapes for 2 different images

3.3. Fusion of the information provided by the two maps

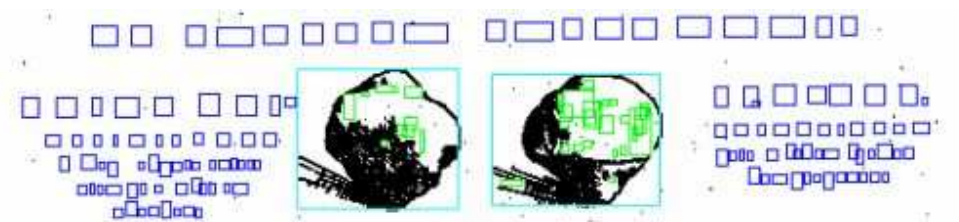
3.3.1. Text/Graphic Separation

The graphical parts do not always correspond to only one connected component in the map of the shapes. On the contrary, most of the time, they correspond to the superposition of the bounding boxes

of several connected components labelled *Text* or *Graphic*. In our processing, when two bounding boxes labelled *Graphic* overlap, they are merged in order to produce only one rectangle including the whole of the detected graphical areas. It is also possible that connected components labelled *Text* are present inside an area labelled *Graphic*. These components are then labelled *Text_Graphic* and would be the subject of a particular analysis to determine if they belong indeed to the *Text* class or to the *Graphic* class. For that, the algorithm that we developed uses once again the concept of map of the background to carry this classification. The algorithm used to distinguish the *Text* components close to or inside a *Graphic* component is as follows:

- 1- *Text* areas inside *Graphic* areas are relabelled *Text_Graphic*
- 2- colouring the *Text* and *Text_Graphic* labelled blocks in white, in the binarised image
- 3- calculation of the background map in the same manner as done previously
- 4- *Text_Graphic* blocks positioned over the dark parts of the map of the background are preserved; the others are removed because they are regarded as pieces of *Graphic* blocks

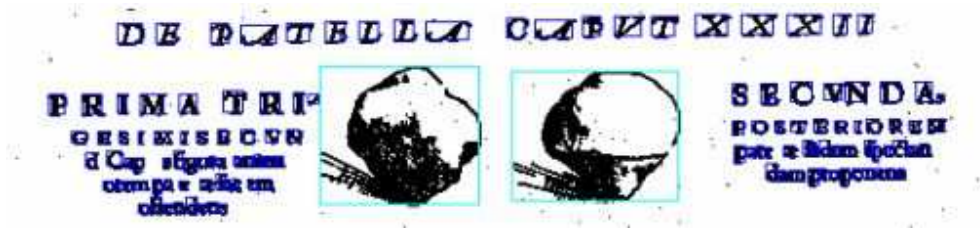
The removed areas are located in strongly textured areas of the binarised image. We assume these areas correspond to the presence of graphical parts in the initial image (see figure 9).



(a) Illustration of the presence of *Text_Graphic* blocks (represented in green)



(b) Map of the background obtained at stage 2



(c) Result: the *Text_Graphic* blocks located inside *Graphic* blocks were removed

Figure 9: Illustration of Text / Graphic separation in graphical parts

3.3.2. Extraction of the text areas

By simultaneously using the information provided by the maps of the shapes and the background, we can extract *Text* areas. We start from the list of the connected components labelled *Text* to rebuild the paragraphs of text by association of connected components likely to be characters. To carry out an association, it is necessary that the two components labelled *Text* are rather close and the segment between their centres of gravity (G_1 and G_2) does not cross an important transition area in the map of the background (low grey level values). This multi-criterion constraint can be expressed by:

$$d(G_1, G_2) \times (256 - \underset{(i, j) \in [G_1, G_2]}{\text{Min}} [Ng(i, j)]) \leq \text{Fusion_Threshold} \quad (2)$$

where d indicates the Euclidean distance between the centres of gravity of the two sets to be associated. When this criterion is checked for two close areas, we give them the same label (attribution of an identical number to each character of the same block of text). If the criterion is not respected, the association is refused. Some other multi-criterion constraints have been tested (such as using the sum of the grey level values instead of the minimum, or the use of the distance between the borders of the 2 areas instead of the centre of gravity, ...) and we kept the best one. The *Fusion_Threshold* can be chosen manually by the users. The selection of this value is quite straightforward. First, the users are advised to start with a low value (which will result in few fusion of the text part) and then later during the interactive analysis, if necessary, the fusion step can be reactivated for difficult text parts. Such strategies are discussed in detail in section 5.2 (Description of the different scenarios).

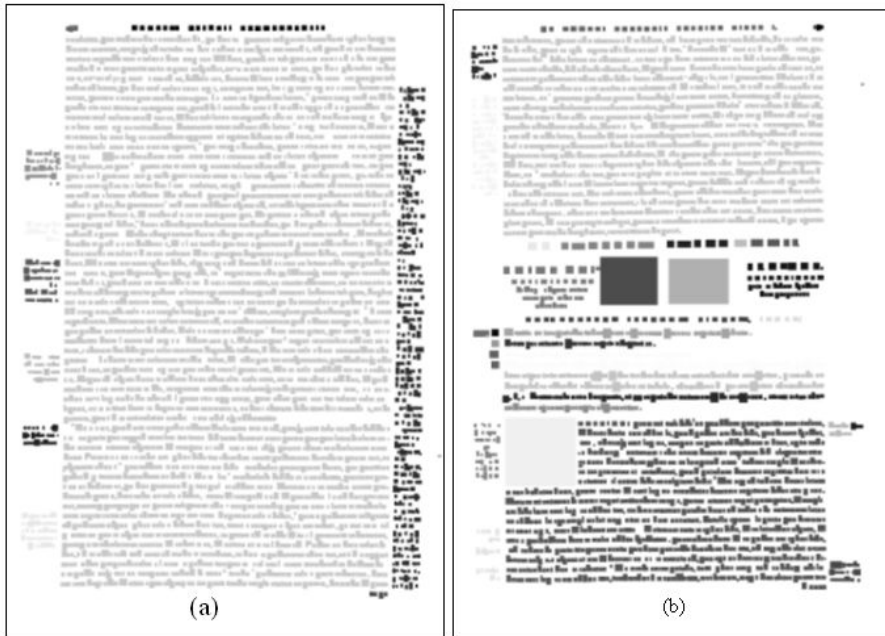


Figure 8: Examples of segmentation results (each segmented block has a different gray level value).

The search of the neighbours between text parts is done successively horizontally (horizontal fusion) then vertically (vertical fusion) in a progressive way and stops when no more fusions are possible. These areas are not always rectangular since they correspond simply to a gathering of the shapes having the same labels. The obtained areas can thus have an arbitrary shape. Two examples of the results are provided in Figure 8. The *Fusion_Threshold* (2) used is the same for the two images. In order to visualize the obtained segmentation, each label given to a component of a given area corresponds to a different grey level value. In example a) all the blocks were correctly segmented in spite of their proximity. In example b), layout of the page and the typesetting rules are much more complex (presence of ornamental letters, images, legends, titles) and our method tends to over-segment the image (there are too many blocks). But, we will see later, this over-segmentation caused by the choice of a strict initial threshold for the fusion is only temporary since the algorithm of fusion can be applied several times with different parameters (insertion of the fusion step in the scenarios of analysis created by the users).

4. User-driven analysis of the layout

Except for extremely specific applications, current block classification or logical layout analysis methods have shown their limitation. Our proposal consists of a general revision of document recognition methods to design a widely usable system (i.e. not being dedicated to a particular model of documents). For this purpose, block classification need to be addressed in a more flexible way. So, we are proposing an architecture that drops the strict processing chain to offer an adequate level of human-machine cooperation. An internal representation of the document structure associated with a set of rules can enable an interactive learning of the model of the documents. Furthermore, the users can use a given example of image to define this new model.

As shown in figure 8, the segmentation stage produces an intermediate representation of the image providing for each page a set of blocks labelled *Text*, *Graphic*, *Text_Graphic* or *Noise*. At this stage of the process, it is possible to ask the user to build different processing sequences (we will name them scenarios of structural analysis) allowing a gradual evolution of this intermediate representation according to the user aims and to the characteristics of the images to analyze. He and Downton have already shown that it is sometime interesting to propose user-assisted document analysis systems [He03]. In the long term, our goal is to obtain automatically the most precise labelling of the contents of the digitized pages of complex documents by the application of a scenario that the user builds interactively and easily (user-driven analysis).

4.1. Interactive setting of scenarios for analysis

Once the initial segmentation of the image has been achieved, the architecture of the system (called AGORA) that we propose, makes it possible to continue the analysis in an interactive way on a typical image. For that, we conceived a set of interfaces that allow the user to build, in his own way, scenarios of analysis. The application of a scenario will allow the progressive evolution (incremental analysis) of the contents of the intermediate representation (initial segmentation) obtained beforehand.

The tools placed at the disposal of the user to build the scenarios are:

- a rule editor allowing the refinement of the labels given to the various obtained areas (*Text*, *Text_Graphic*, *Graphic*, *Noise*). These rules are observed in a sequential way according to a strategy defined by the user who can, for example:
 - be interested by the *Graphical parts* only or on the contrary by the *Text* parts only
 - choose to label first the easily characterisable blocks using simple rules. Then, he can use a more complete context to extract the less easily identifiable parts of the document [Ramel98].
- the application of the algorithms for horizontal or vertical fusion with various thresholds and only on a specific type of blocks defined by the user (for example to merge only the blocks labelled *Title* which were split at the beginning of the process in several blocks as in the example of figure 8)
- deletion of the blocks of a particular type (for example the *Text* labelled blocks if the user is interested only in the ornamental letters)

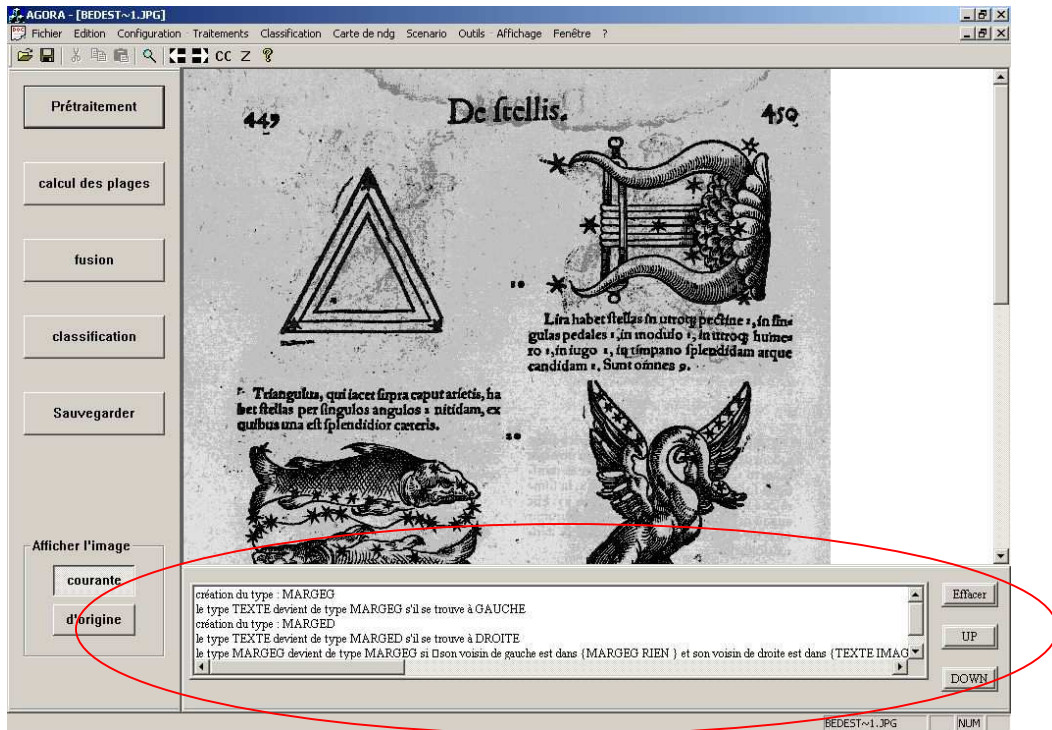


Figure 10: View of the Agora software. The rules of the active scenario is displayed

at the bottom of the window

To build a scenario, the user performs (as for a macro) the successive actions which have to be recorded on a typical image. These actions are applied to the image and the results are displayed in real time. We can thus validate their effectiveness. The actions (rules) are translated in a literal list permanently displayed (see figure 10). The user can, then, manage this list to modify the scenario by reordering or removing processing rules. Once the scenario is considered as correct, it can be saved in a file for storage or for application on a more subsequent set of images (a complete book) during a batch processing.

Currently the usable rules to make the labels of the blocks contained in the intermediate representation concern (figure 11):

- **the topological position** of the blocks
- **the neighbourhood relations** between identified blocks
- **the shape and the content** of the blocks

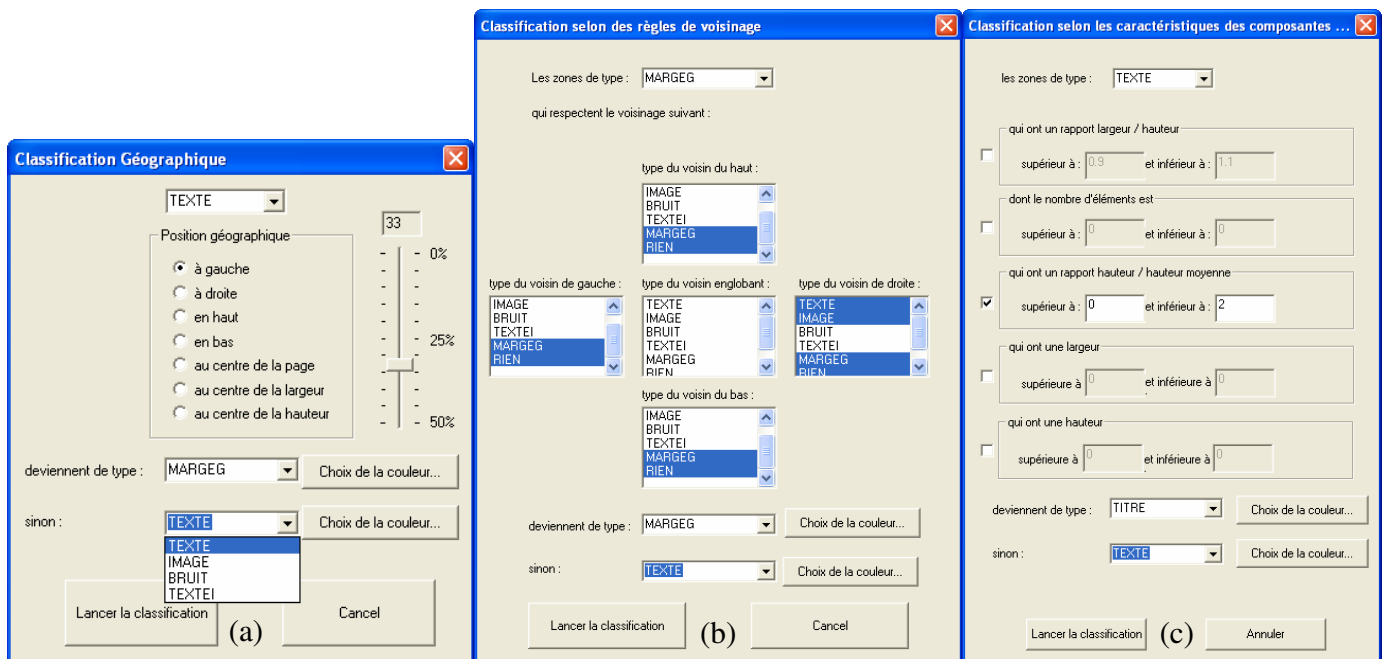


Figure 11: Interfaces for the installation of position (a), neighbourhood (b), and intrinsic properties (c) rules

4.2. Study of the topological position

It is possible to identify a number of invariants on the topological positioning of the objects present in a historical book making it possible to associate a label to them. The goal is not to use extremely strict rules because page layout can be variable (it is unwise to specify that the centre of gravity of a note in left margin is between the pixel of X-coordinate 205 and 213). However, one can reasonably specify that the position of the centre of gravity of a left margin text is doubtless located in the first third of the width of the page. The suggested interface (see figure 11a) makes it possible to take account of the topological position (left, right, top, bottom, centre) of the centre of gravity of the blocks to make the label evolve. Of course this type of rule will only constitute a first index (provisional label) to correctly isolate an object of a given type.

4.3. Study of the neighbourhood

It is also possible to insert rules in a scenario which concern the neighbourhood relations between blocks. Examples of rules allowing to label and extract different elements are provided in figure 12. For example, the second interface (figure 11b) allows to define rules like:

A *Graphic* block will be labelled *Ornamental letter* if it has :

- at the left, a block labelled *Left Margin* OR *Column* OR *Graphic* OR *Nothing*
- AND at the right, a block labelled *Legend* OR *Column* OR *Graphic*
- AND above, a block labelled *Legend* OR *Column* OR *Title* OR *Text* OR *Graphic*
- AND below, a block labelled *Legend* OR *Column* OR *Graphic*

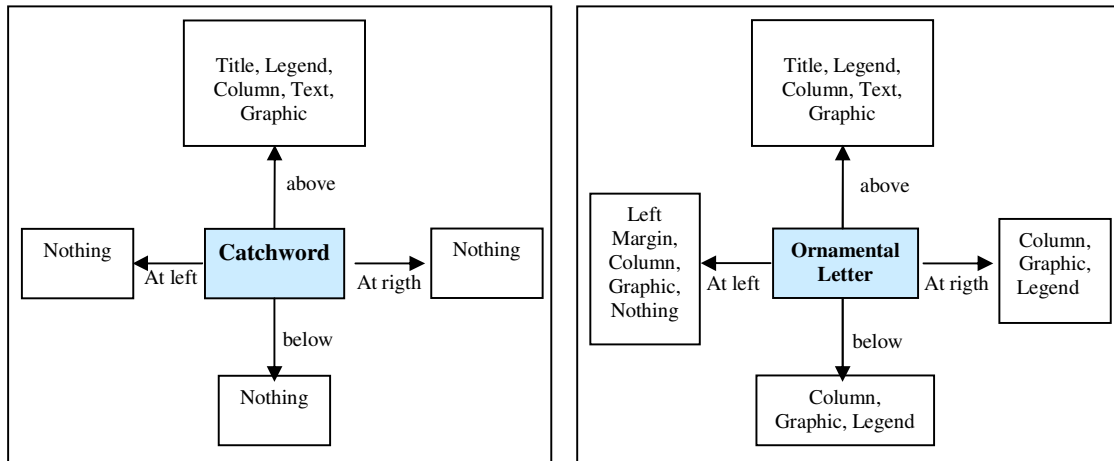


Figure 12: Examples of neighbourhood rules

4.4. Study of shapes and contents

Another interface (figure 11c) makes it possible to use intrinsic properties of the shape to refine the labelling. The actual and non exhaustive proposed criteria are

- the *size* of the block (*width* and *height*),
- the *Height / Width ratio* of a block, for example for the location of the square blocks (ornamental letters) or of the blocks having a width larger than height (border)
- the *number of connected components* (elements) counted in the block that can prove to be interesting in identifying certain objects
- the ratio of *block Height / average Height of the connected components* can be usually effective to locate the textual blocks comprising only one line of text (titles, legends...).

To finish this part, let us recall that, in addition to these rules allowing evolution of the labels of blocks, it is possible to insert rules of fusion and suppression of blocks in the scenarios. Examples of scenarios built by users, maybe not specialists in image processing can be seen in the following part.

5. Experimentations and results

Our software was made available to the CESR which currently uses it in an intensive way to process, analyse, index and make available on line their historical books. Training in use of the software

(segmentation and creation of scenarios) was given to the potential users (historians, librarians, keyboard operators ...) so that they could produce scenarios and consequently test our user-driven system on numerous images. This collaboration between the CESR and our laboratory makes it possible to improve and complete the interfaces of the software taking into account the needs of the final users.

Several experiments were undertaken, either by the staff of the CESR, or within our laboratory. The first experiment was devoted to the analysis of the robustness of the block segmentation with respect to the skew. This test has been achieved using 10 images of current documents and 10 images of historical books on which we have applied rotations from 1 to 20 degrees. A second experiment has been carried out in our lab on a complete historical book with 250 pages. The objective of this test was to validate the usability of the proposed interfaces to build different scenarios of analysis for text / graphic separation and classification.

AGORA has then been provided to the real end-users at the CESR where numerous experiments have been done. The results of some of them are presented in this article: the first one aimed at evaluating whether only one scenario could be used to analyse different books or it is preferable to settle one scenario for each book. To conclude this part, two different scenarios are commented to illustrate the flexibility of the proposed system. All these experiments have led to the results presented in the next sections.

5.1. Initial segmentation

We call "initial segmentation", the result obtained after the first application of the algorithm of fusion on the connected components according to the foreground and the background maps which provides a set of blocks labelled *Text*, *Graphic*, *Text_Graphic* and *Noise*. This segmentation is not the final one (since it will evolve during the scenarios of incremental analysis). So, we evaluate, in this part, only the resistance of the first step to the skew problems. Next, the results obtained after the complete application of a scenario will prove the global performance of the AGORA system.

We tested the resistance of our segmentation algorithm to the skew of images on old documents and on images of current periodicals. The results obtained using 20 images are summarized in figures 13 and 14 and indicate a sufficient tolerance for the old documents and a strong tolerance for more structured documents. If the documents are badly positioned at the time of scanning or if the images are badly illuminated, it is preferable to apply a correction of the geometrical defects before using the segmentation algorithm. When the image scans are correct, the obtained results are very satisfactory without pre-processing.

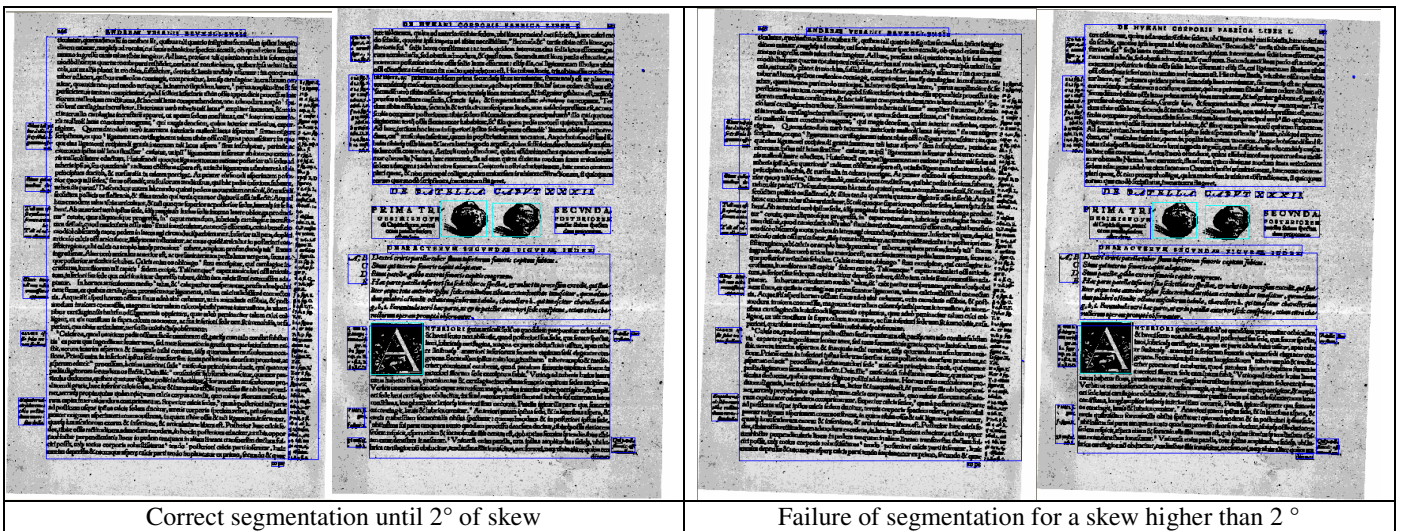


Figure 13: Sensitivity to rotation on images of old documents

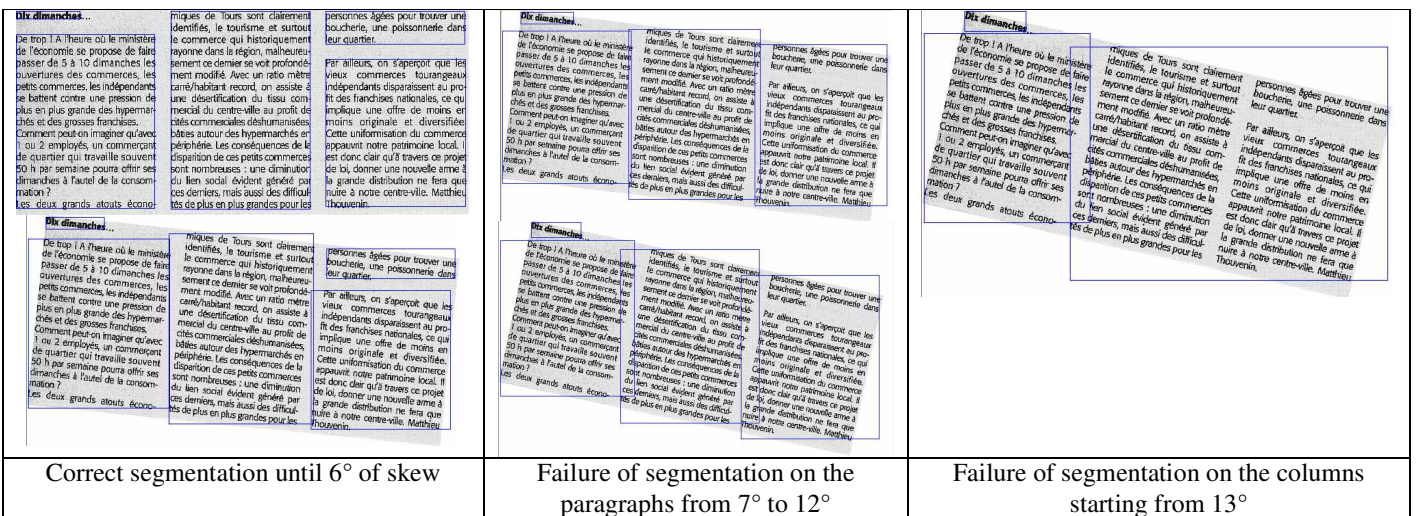


Figure 14: Sensitivity to skew on images of current documents

5.2. User-driven structural analysis

During a first experiment carried out in our laboratory, a scenario allowing the labelling of seven types of historical elements was set up : " Left Margin (MG)", " Right Margin (MD)", " Ornamental letter (Ltr)", " Right Page Number (NPD)", " Left Page Number (NPG)", " Principal Title (TP)" and " Title (T) ". For this test of block classification, we use an image database of 250 pages coming from a complex historical book provided by the CESR. Previous examples of pages used to illustrate this article have been extracted from this book. The obtained results are provided in figure 15. The tiresome work of checking the assigned class for each zone was carried out in a visual way on each analyzed page. Table 1 shows the results obtained during the batch processing which lasted 9h20mins on a Athlon barton 2.5 Ghz with 768 MB of RAM

Detected Types	Detected	Not detected	False detections	Rate of detection (%)
Margins	207	18	0	92
Numbers of right page	97	32	7	74
Numbers of left page	106	30	10	76
Principal titles	223	27	0	90
Ornamental letters	80	1	0	99

Table 1: Results obtained with scenario 1

These results are dependant on the initial segmentation and on the robustness of the rules constituting the scenario built by the users. The aim of the user who conceived this scenario was not to have any bad classification (minimal risk) for the *Ornamental Letters* and the *Margins*.

The scenario used was thoughtful but did not implement "redundant rules" which could have involved additional risk on certain criteria. These "redundant rules" aids in the verification of the high level elements (ornamental letters, for example) that have already been extracted with the help of the processing rules. For example once ornamental letters have been located using rules describing their size and positions in the image, user can specify a redundant rule to verify that no ornamental letter

have been found to the right of another ornamental letter(impossible in real images). If found, the redundant rule will change the label “ornamental letter” with the label “unknown graphic”, for example. Indeed, the labels generated at a given moment can be validated later using other rules. With this scenario, 92% of *the Margins* are extracted correctly without any wrong detection. This scenario could be enhanced in order to locate the missing *Margins* once again by adding more rules when the other high level objects would have been localised (evolution of the context to facilitate new extractions). Still with the same scenario, 99% of *the Ornamental letters* were well classified. This is satisfactory since the labelling of *the Ornamental letters* uses the result of the classification of *the Margins*. The table 1 shows that the *Page Numbers* were difficult to extract or that the rules used to locate them were not sufficient.

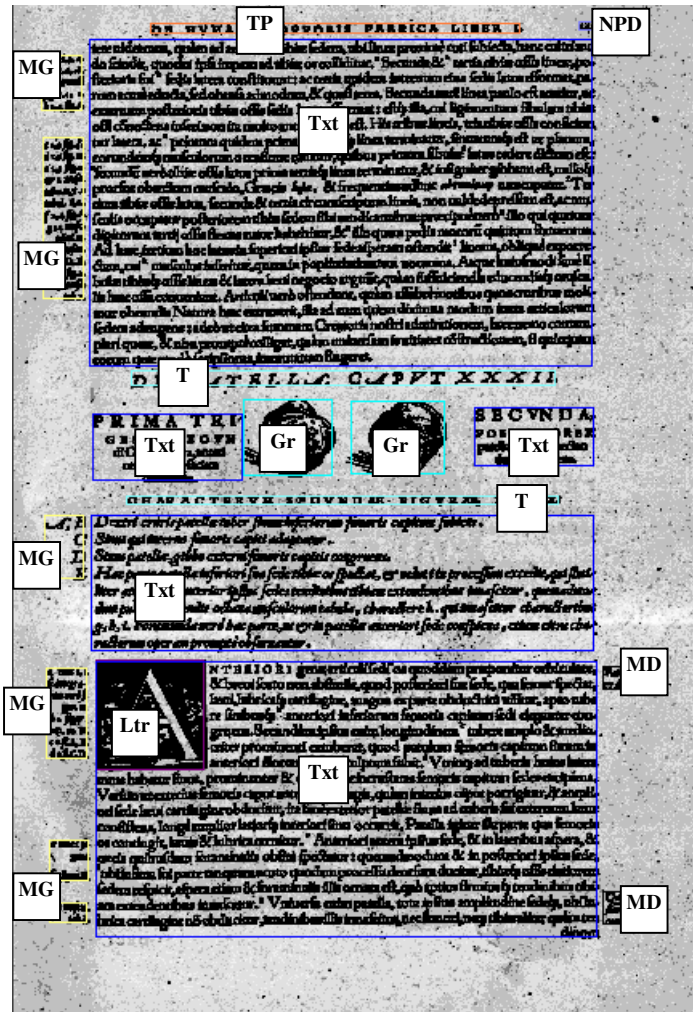


Figure 15: Example of result obtained with a first scenario 1 (Left Margin (MG), Right Margin (MD), Ornamental letter (Ltr), Right Page Number (NPD), Left Page Number (NPG), Principal Title (TP) and Title (T)).

Many other experiments were carried out at the CESR by the staff following the training in the use of the software. In addition to the recognition rates, it is interesting to notice the way in which the produced scenarios are structured and applied.

At the CESR, all the tests were carried out on images resulting from a digital camera that provides 1200 X 2000 grey level pixels images. The segmentation parameters were as follows:

- automatic binarisation
- minimum height of a large component: 60
- minimum width of a large component: 60
- maximum height of a small component: 5
- maximum width of a small component: 5
- Threshold of horizontal fusion: 500
- Threshold of vertical fusion: 500

The first result relates to the initial segmentation which did not produce any error of labelling: All the *Text* and *Graphic* blocks were correctly detected in all the books selected, however, sometimes an over-segmentation of some blocks resulted from the use of the strict fusion thresholds (500). But, one can notice that more fusions can be achieved later during scenarios of classification to correct this over-segmentation.

First scenario

A first scenario was carried out on 1452 images from five different books. The objective of this scenario (shown below) was to identify the various types of the graphical blocks present in the five different books. As explained in section 3.1, at the end of the first step of the processing, the representations of the image contain a list of blocks labelled as Text, Noise or Graphic. Then, the user has to settle scenarios to make these initial labels evolving using specific rules as shown below:

-Text on the left of the page with 15% becomes Left Margin (15% of the page width is a threshold defined by the user in this scenario)

-Text on the right of the page with 15% becomes Right Margin

- Vertical Fusion of the Left and of the Right Margins with a threshold of 100000

- Graphic with at the Left : Left Margin OR Nothing AND at the right : Text AND Width/Height Ratio between $0,8 < < 1,2$ becomes Ornamental Letter

- Graphic with Width/Height Ratio between $3 < < 10$ become Border

- Graphic in the centre of the page become Portrait

- Portrait with above : Portrait OR below : Portrait OR at the bottom of the page with 20% becomes Graphic

- Graphic with at the Left : Text OR at the right : Text OR above : Text OR below : Text OR Nothing AND in the centre of the page becomes Floret

- Suppression of the Text and the Left and Right Margins blocks

Here, we can notice that the first rules often use only topological properties of the blocks to change their labels. This probably means that these kinds of rules are easy to understand and manage (but perhaps not very efficient). When *Margins* are detected, it is possible to run for another time the fusion step (used during the segmentation) to merge only *Margin* blocks. In this scenario neighbourhood and shapes rules are used to locate *Ornamental Letters* and *Borders*. The images were processed in 11h22 and the results obtained are shown in table 2.





Type		Detected	False Detections	Not Detected	Rate of detection
Border		81	4	8	90.6 %
Floret		36	2	0	100 %
Ornamental letter		294	56	11	95.6 %
Portrait		89	31	0	100 %

Table 2: Results obtained on the 1452 pages of the five books. The books were mainly textbooks and containing a variety of 85 borders, 34 florets, 249 ornamental letters and 58 portraits

Taking into account the diversity of the selected old books (size, layout...) and the simplicity of this scenario, the results appear quite good. Nevertheless, to obtain better detection rates and to extract more ambiguous objects, it is often preferable to adapt the scenario to each book. Thus, the continuation of the tests was achieved on images coming from a unique book. Some pages of this book are presented in figure 16.

Second scenario



Figure 16: Example of pages from the selected book after the pre-processing with Book Restorer

A second test was carried out on 180 pages of the selected book. Because some images are bad illuminated or some text lines have an important curvature, Book Restorer [Digibook2005]) has been used for the pre-processing of the images (lighting correction for the spots and the shade of the binding, geometrical correction of the curvature):

- Graphic with Width/Height Ratio between $3 < < 10$ becomes Border
- Graphic with at the left : Left Margin OR Nothing AND at the Right : Text AND Width/Height Ratio between $0,7 < < 1,3$ becomes Ornamental Letter
- Vertical Fusion of Text with a threshold of 2000
- Text on the left of the page with 25% becomes Left Margin - Text on the right of the page with 25% becomes Right Margin
- Vertical Fusion of the Margins with a threshold of 100000
- Horizontal Fusion of the Text with a threshold of 2000
- Text on top of the page with 10% AND with a number of elements between $0 < < 4$ becomes Pagination
- Text on bottom of the page with 25% AND with a number of elements between $0 < < 6$ becomes Signature
- Signature with Text below, on the left or on the right AND on the left of the page with 50% becomes Text
- Vertical Fusion of the Text with a threshold of 3000
- Horizontal Fusion of the Text with a threshold of 3000
- Text with height in pixel between $0 < < 100$ AND on the top of the page with 15% becomes Title
- Horizontal Fusion of the Title with a threshold of 3000
- Suppression of the Title if the number of elements is lower than 8
- Suppression of the Title if the Width/Height Ratio is lower than 1
- Suppression of the Title if it is on the left of the page with 10%

- *Suppression of the Text*

This scenario is again time quite simple. All kinds of rules (topological, shape, neighbourhood) are equally used. Fusion steps are used very often to merge specific block like *Margin* and *Text*. The processing lasted 2h32mins and produced the following results:

- Detection of the *Borders* with 100%
- Detection of the *Ornamental Letters* with 96%
- Detected *Pagination* with 79%
- Identification of 100% of the *Titles* but with 4% of noise elements (spots, noises...)

Third scenario

A last test was applied to the 1202 images without pre-processing, with the objective of only labelling the *Text* blocks. The following scenario was used:

- *Suppression of the Graphics*

- *Horizontal Fusion of the Text with a threshold of 3000*

- *Text in the left of the page with 25% becomes Left Margin - Text on the right of the page with 25% becomes Right Margin*

- *Vertical Fusion of the Margins with a threshold of 3000*

- *Margins with a Width/Height Ratio between $4 < < 30$ becomes Text*

- *Text on bottom of the page with 25% AND with the number of elements between $0 < < 8$ becomes Signature*

- *Vertical Fusion of the Text with a threshold of 3000*

- *Text with Width/Height average in $0 < < 2$ AND on top of the page with 15% AND number of elements is lower than 50 becomes Title*

- *Horizontal Fusion of the type Titles with a threshold of 3000*

After 10h52mins of processing time, all the *Titles* and *Signatures* were detected correctly. The extraction of the *Margins* functioned but produced many spurious elements (shades, spots...)

Discussion

To conclude this experimental stage, we can first notice that the average processing time for one image is 30 seconds. The obtained results are better on pre-processed images. The automatic binarisation

(by Sauvola algorithm [Sauvola00]) increases the processing time by 25 to 35 % as well as increasing the precision of the detection. Lastly, it is preferable not to merge the blocks too much at the beginning (initial segmentation) in order to be able, in certain cases, to isolate small elements (*Signature* or *Pagination* for example) before carrying out additional fusions later in the scenario of analysis. The presented scenarios show that the users prefer to implement several scenarios (one for the text part and one for the graphical part) instead of using only one global scenario even if it would probably have provided better results. It probably means that it is not as easy to build thoughtful scenarios as one may initially think.

Seeing these experiments and results, we think the main advantages coming from the use of AGORA are :

- Robustness to proximity between blocks during the segmentation step. Our hybrid method of segmentation avoids wrong fusion of close text blocks
- Working even if blocks are not rectangular and if the layout is variable
- Robustness to normal skew (for a skew higher than 2° , a skew correction step is needed)
- Powerful separation of Text/ Non-text areas even for characters in graphical parts
- Of course, genericity due to the use of a “user-driven method”. With the help of scenarios, users can label every type of components (ornamental letters or catchwords in historical books, titles and paragraphs in current books) and not only paragraphs, tables and graphics like with other systems [Dori97, Hadjar01, Marinai05, Wang06].

The only weak point of AGORA shown by our experiments comes from the binarisation step that can be sensitive to bad illumination or to noise coming from the degradations in the initial documents.

6. Conclusion

In the first part of this article, we highlighted the sources of errors in the traditional methods of page decomposition using a characterization of page layout in the historical books. We noticed that each type of method (bottom-up and top-down) required different information which one should have been

aware of to reach a segmentation of quality. For that, our recommendation consists of using first a hybrid algorithm based on the construction of two representations of the image: the map of the shapes which is focused on the connected components and the background map which provides information on white spaces separating the blocks constituting the page. The joint analysis of the contents of these two maps makes it possible to lead to a robust initial segmentation of the image. The results obtained with this method are very interesting; the adjustment of the necessary parameters is straightforward and not sensitive to variations.

Second, the originality of our approach lies in the opportunity which we offer to the users to be able to build, in an interactive way, scenarios of incremental analysis. We propose to call this new method “user-driven analysis” in opposition to data-driven or model-driven methods. The goal is, on the basis of the initial segmentation, to be able to make the representation of the images evolve in a progressive way to lead to the finest possible characterization of its contents according to the user objectives and to the type of images to be analyzed. The CESR has processed several complete books using AGORA prototype and their own scenarios of block classification. Thus, the CESR has increased the number of books offered to the users in its Virtual Library (see <http://www.bvh.univ-tours.fr>). Even if the system produced some errors, the processing and time saved as compared to manual processing is considerable (for example, the manual indexation of the page layout of an historical book of 300 pages last approximately two days instead of only two hours when using Agora), this providing to the specialists of historical books, a useful tool which they had never imagine.

7. References

[Akindele93] O. Akindele, A. Belaid. *Page Segmentation by Segment Tracing* In Proc. of the 2nd International Conference on Document Analysis and Recognition, p341-344, 1993.

[Antona98] A. Antonacopoulos, *Page Segmentation Using the Description of the Background*. Computer Vision and Image Understanding, Special Issue on Document Image Understanding and Retrieval, Vol. 70, No. 3, p350-369, 1998.

- [Baird90]. H. S. Baird, S. E. Jones, and S. J. Fortune, *Image segmentation by shape-directed covers*. In Proceedings of International Conference on Pattern Recognition, p820-825, (Atlantic City, NJ), June 1990.
- [Baird92] H Baird. *Background structure in document images*. In Advances in Structural and Syntactical Pattern Recognition, ED. H. Bunke. p253-269. 1992.
- [Belaid97] A. Belaïd, *Computer aided design of models of page for their use in recognition of documents*, Workshop one Electronic Page Models (LAMPE' 97). 1997.
- [Breuel03] T. M. Breuel. *An Algorithm for Finding Maximal Whitespace Rectangles at Arbitrary Orientations for Document Layout Analysis*. International Conference for Document Analysis and Recognition (ICDAR 2003). Vol.1, p66-70. 2003.
- [Chen95] S. Chen and R.M. Haralick. Recursive Erosion, Dilation, Opening and Closing Transforms. IEEE Transaction on Image Processing, Vol 4(3):p335-345, 1995.
- [Digibook2005] Book scanner web site : <http://www.i2s-bookscanner.com/fr/default.asp>
- [Dori97] D. Dori, D. Doermann, C. Shin, R. Haralick, I. Phillips, M. Buchman, and D. Ross. *The Representation of Document Structure: A Generic Object-Process Analysis*. In H. Bunke and P. S. P. Wang, editors, Handbook of Character Recognition and Document Image Analysis, chapter 16, pages 421--456. World Scientific, 1997.
- [Eternad97] Kamran Etemad, David Doermann, Rama Chellappa, *Multiscale Segmentation of Unstructured Document Pages Using Soft Decision Integration*, IEEE Transactions on Pattern Analysis and Machine Intelligence, v.19 n.1, p92-96, January 1997
- [Jain92] A.K. Jain and S. Bhattacharjee. *Text Segmentation using Gabor fillters for automatic document processing*. Machine Vision and Applications, Vol. 5(3): p169-184, 1992.
- [Jain96] Jain, A.K., Zhong, Y., *Page Segmentation Using Texture Analysis*, Pattern Recognition. Vol. 29, No. 5, p743-770, 1996.
- [Hadjjar01] K Hadjar, O Hitz, R. Ingold. *Newspaper page decomposition using Split and merge approach*. Proceedings of the 5th International Conference one Document Analysis and Recognition. p1186-1191, 2001.
- [Hadjjar02] K Hadjar, O Hitz, L Robadey, R. Ingold. *Configuration REcognition Model for Complex Transfers Methods Engineering: 2(CREM)*. Proceedings of the 5th International Workshop one Document Analysis Systems. p469-479, 2002 .

- [He03] J. He, A. Downton: *User-Assisted Archive Document Image Analysis for Digital Library Construction*. Proceedings of the 6th International Conference on Document Analysis and Recognition.. p498-502, 2003:
- [Journet2005] N. Journet, V. Eglin, J.Y. Ramel, R. Mullot. *Text/Graphic labelling of Ancient Printed Documents*. International Conference on Document Analysis and Recognition (ICDAR'05). Séoul. Corée. p1010-1014, 2005.
- [Kise98] K. Kise, A. Sato, M. Iwata. *Segmentation of page images using the area Voronoi diagram* Computer Vision and Image Understanding Special issue on document image understanding and retrieval. Vol. 70(3). p370–382, 1998.
- [Le94] D.S. Le, G.R. Thoma, and H. Wechsler. *Automated Page Orientation and Skew Angle detection for Binary Document Images*. Pattern Recognition, 27(10):p1325-1344, 1994.
- [Lebourgeois99] F. LeBourgeois, H. Emptoz, *Document Analysis in Gray level and typography extraction using Character Pattern redundancies*, In proc. of the 5th International Conference on Document Analysis and Recognition., Bangalore India, p177-180. 1999.
- [Lebourgeois03] F. Lebourgeois, H. Emptoz, E. Trinh, *Compression et accessibilité aux images de documents numérisés / Application au projet Debora*. Document Numérique. Vol 7(3-4). p103-127. 2003
- [Lee01] Lee, S.W., Ryu, D.S., *Parameter-Free Geometric Document Layout Analysis*, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 23, No. 11, p1240-1256, 2001.
- [Marinai05] Simone Marinai, Marco Gori, Giovanni Soda: *Artificial Neural Networks for Document Analysis and Recognition*. IEEE Trans. Pattern Analysis and Machine Intelligence. Vol.27(1): p23-35 2005.
- [Min96] Y. Min, S.-B. Cho, and Y. Lee. *A Data Reduction Method for Efficient Document Skew Estimation Based on Hough Transformation*. In Proc. of the 13th International Conference on Pattern Recognition, Vienna, Austria, p732-736, 1996.
- [Nagy84] G. Nagy and S. Seth. *Hierarchical representation of optically scanned documents*. In 7th International Conference on Pattern Recognition (ICPR), p347-349, 1984.
- [Nagy93] G. Nagy, S. Seth, M. Krishnamoorthy, and M. Viswanathan. *Syntactic Segmentation and Labeling of Digitized Pages from Technical Journals*. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 15(7): p737-747, 1993.
- [Ogorman93] L O' Gorman. *The Document Spectrum for Page Analysis Layout* IEEE Transactions

on Pattern Analysis and Machine Intelligence, Vol. 15(11). p1162-1173. 1993

[OGorman95] L. O'Gorman and R. Kasturi. *Document Image Analysis*. IEEE Computer Society Press, Los Alamitos, CA, 1995.

[Ramel98] JY Ramel, N. Vincent, H. Emptoz *Extraction contextuelle d'entités graphiques dans les dessins : du plus simple au plus complexe....* Colloque International Francophone sur l'Écrit et le Document. Quebec (Canada). p453-462. 1998.

[Trinh03] E. Trinh, *De la numérisation à la consultation de documents anciens*. Thèse de doctorat en Informatique. Insa de Lyon. 2003

[Sauvola00] Sauvola, J., Pietikainen, M.: *Adaptive Document Image Binarisation*. Pattern Recognition Vol. 33, p225-236. 2000

[Spitz98] Spitz A.L., *Analysis of Compressed Document Images for Dominant Skew, Multiple Skew and Logotype Detection*, Computer Vision and Image Understanding Vol. 70, 3, p321-334, 1998.

[Wang89] D. Wang and S.N. Srihari. *Classification of Newspaper Image Blocks Using Texture Analysis*. Computer Vision, Graphics and Image Processing, Vol. 47, p327-352, 1989.

[Wang95] S.-Y. Wang and T. Yagasaki. *Block Selection: A Method for Segmenting Page Image of Various Editing Styles*. In Proc. of the 3th International Conference on Document Analysis and Recognition, Montreal, Canada, p128-133, 1995.

[Wang06] Y. Wang, IT. Phillips, RM. Haralick, *Document zone content classification and its performance evaluation*. Pattern Recognition. Vol 39. p57-73. 2006.

[Wong82] K.Y. Wong, R.G. Casey, and F.M. Wahl, *Document Analysis System*, IBM Journal of Research and Development, vol. 26, no. 6, p647-656, Nov. 1982.