



**Ancient document image segmentation using
the autocorrelation function and multiresolution analysis**

Maroua MEHRI*†

Supervised by Petra GOMEZ-KRÄMER*, Pierre HÉROUX†,
Alain BOUCHER*, and Rémy MULLOT*

Funded by **ANR-DIGIDOC project**

*. L3i, University of La Rochelle, La Rochelle, France

†. LITIS, University of Rouen, Saint-Etienne-du-Rouvray, France

E-mails: {maroua.mehri, petra.gomez, alain.boucher, remy.mullot}@univ-lr.fr, and pierre.heroux@univ-rouen.fr

Objective

DIGIDOC : Document Image diGitisation with Interactive
DescriptiOn Capability

Objective

DIGIDOC : Document Image diGitisation with Interactive DescriptiOn Capability

- Control the quality of historical document image digitization

Objective

DIGIDOC : Document Image diGitisation with Interactive DescriptiOn Capability

- Control the quality of historical document image digitization
- Construct a computer-aided categorization tool of pages

Objective

DIGIDOC : Document Image diGitisation with Interactive DescriptiOn Capability

- Control the quality of historical document image digitization
- Construct a computer-aided categorization tool of pages
- Provide a similarity measure between pages

Objective

DIGIDOC : Document Image diGitisation with Interactive DescriptiOn Capability

- Control the quality of historical document image digitization
- Construct a computer-aided categorization tool of pages
- Provide a similarity measure between pages
- **Characterize the document content using intermediate level metadata**

Overview

Book



Overview

Book



- Without any information about the document structure (model)
- Without knowledge of typographical parameters (font size)

Overview

Book



- Without any information about the document structure (model)
- Without knowledge of typographical parameters (font size)



Overview

Book

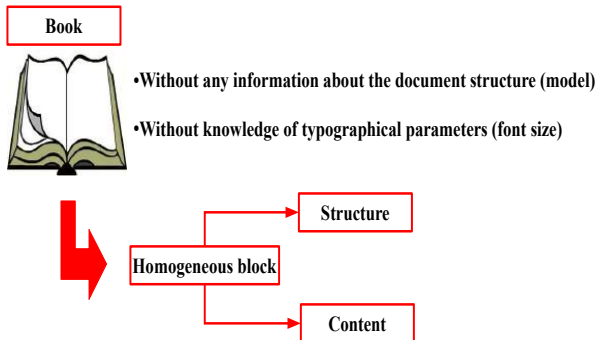


- Without any information about the document structure (model)
- Without knowledge of typographical parameters (font size)



Homogeneous block

Overview

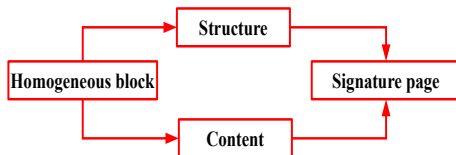


Overview

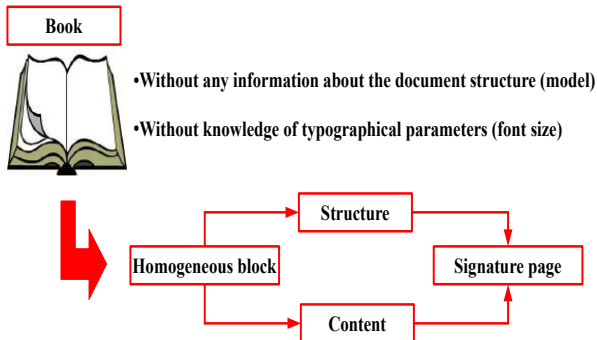
Book



- Without any information about the document structure (model)
- Without knowledge of typographical parameters (font size)



Overview



⇒ Our goal is to propose a set of metadata characterizing the **physical structure of pages** in terms of **homogeneous blocks** and **topological relationships**

Overview

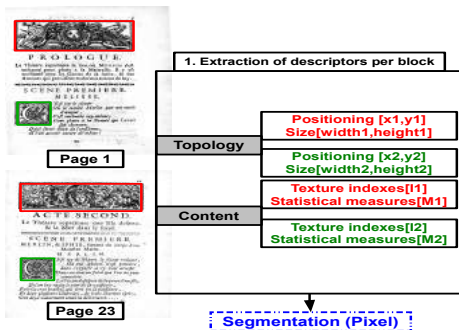


Page 1

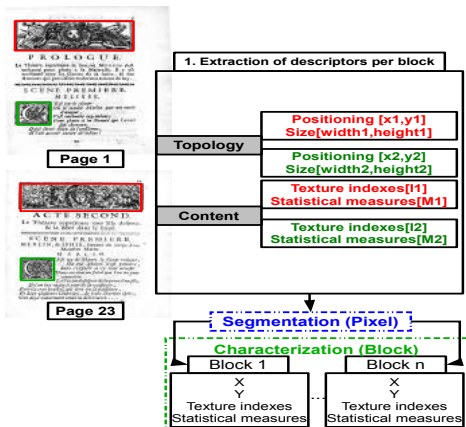


Page 23

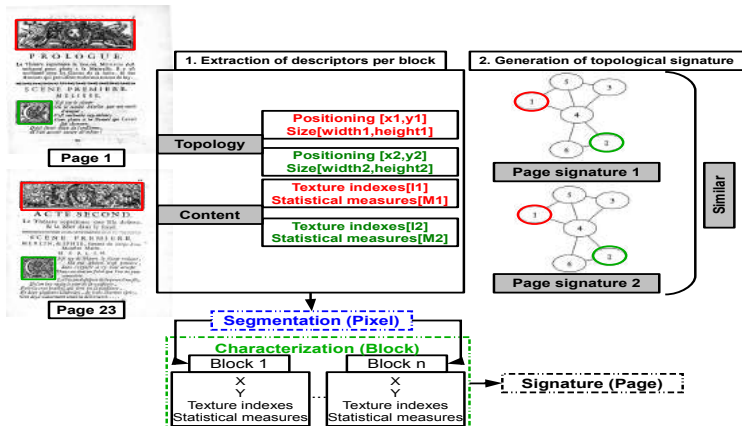
Overview



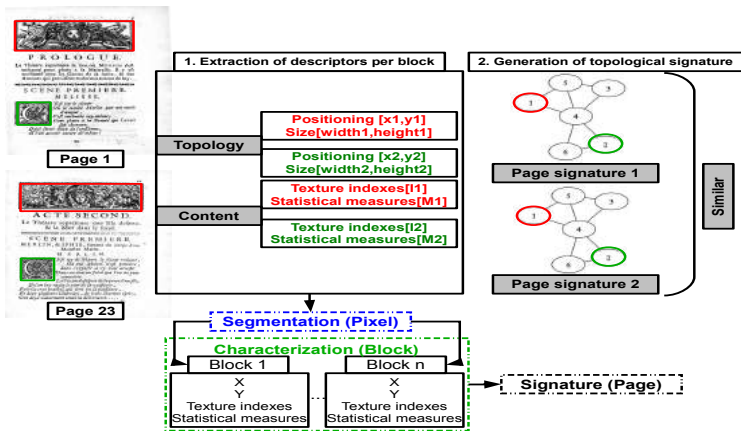
Overview



Overview

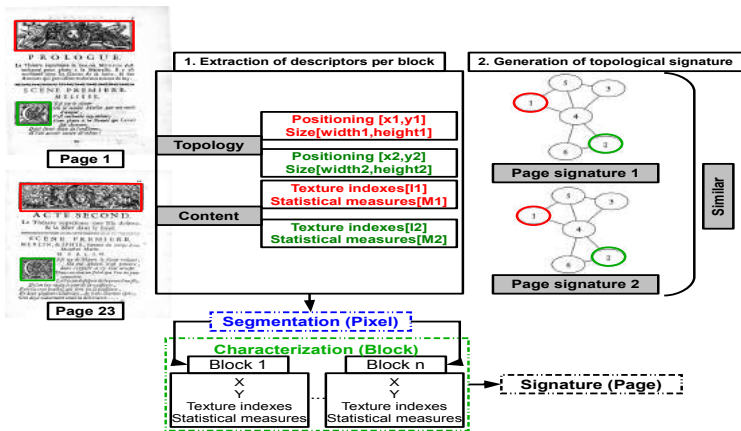


Overview



1 Extraction of descriptors per block

Overview



- 1 Extraction of descriptors per block
- 2 Generation of topological signature

Outline

1 Method

Outline

1 Method

2 Evaluation

Outline

- 1 **Method**
- 2 **Evaluation**
- 3 **Conclusion and further work**

Outline

1 Method

- Previous work
- Proposed method
 - Segmentation using the autocorrelation function and multiresolution
 - Unsupervised clustering approach

2 Evaluation

- Corpus
 - Overview
 - Examples
- Experimental protocol
 - One font and graphics
 - Two fonts and graphics
 - Only two fonts
- Evaluation
 - Homogeneity measure
 - Example of segmentation and evaluation result
- Results
 - Examples of segmentation result
 - Accuracy metrics

3 Conclusion and further work

- Conclusion and further work
- Contribution

Previous work

- Traditional segmentation methods

Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]

Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]

Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]
 - Split-and-merge algorithm [9]

Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]
 - Split-and-merge algorithm [9]
- Segmentation methods based on texture analysis

Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]
 - Split-and-merge algorithm [9]
- Segmentation methods based on texture analysis
 - Statistical (GLCM (Grey Level Co-occurrence Matrix) [3])

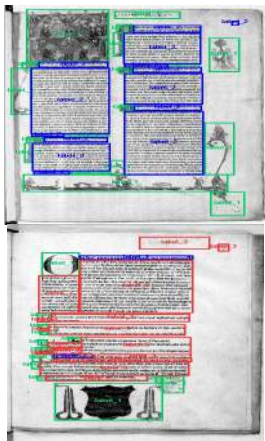
Previous work

- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]
 - Split-and-merge algorithm [9]
- Segmentation methods based on texture analysis
 - Statistical (GLCM (Grey Level Co-occurrence Matrix) [3])
 - Geometrical (Difference-of-Gaussian filter [13])
 - Model-based (Markov random fields [7] and fractals [2])

Previous work

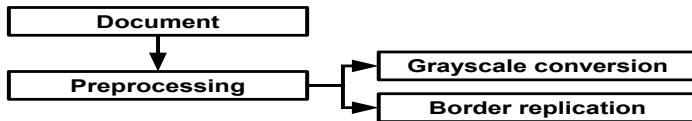
- Traditional segmentation methods
 - Run Length Smearing Algorithm (RLSA) [14]
 - XY-CUT algorithm [6]
 - Split-and-merge algorithm [9]
- Segmentation methods based on texture analysis
 - Statistical (GLCM (Grey Level Co-occurrence Matrix) [3])
 - Geometrical (Difference-of-Gaussian filter [13])
 - Model-based (Markov random fields [7] and fractals [2])
 - Signal processing (Autocorrelation [5], Gabor filters [4], Fourier transforms [10], wavelets [10], and moment-based texture segmentation [12])

Segmentation methods based on texture analysis

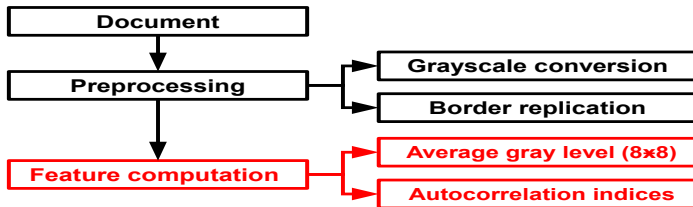


[5] N. Journet, J. Ramel, R. Mullot, and V. Eglin, “Document image characterization using a multiresolution analysis of the texture: application to old documents,” in *International Journal of Document Analysis and Recognition (IJ DAR 2008)*. Springer-Verlag, 2008, pp. 9–18

Autocorrelation feature extraction



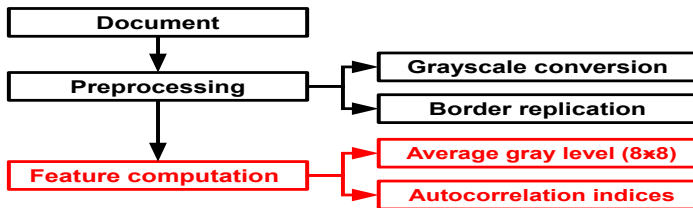
Autocorrelation feature extraction



Feature vector

- Main Orientation: $F^{(1)}_{(x,y)}$
- Autocorrelation intensity: $F^{(2)}_{(x,y)}$
- Directional rose variance: $F^{(3)}_{(x,y)}$
- Mean stroke width: $F^{(4)}_{(x,y)}$
- Mean stroke height: $F^{(5)}_{(x,y)}$

Autocorrelation feature extraction



Feature vector

- Main Orientation: $F^{(1)}(x,y)$
- Autocorrelation intensity: $F^{(2)}(x,y)$
- Directional rose variance: $F^{(3)}(x,y)$
- Mean stroke width: $F^{(4)}(x,y)$
- Mean stroke height: $F^{(5)}(x,y)$

Sliding Windows

- (16x16)
- (32x32)
- (64x64)
- (128x128)

Autocorrelation function & Directional rose

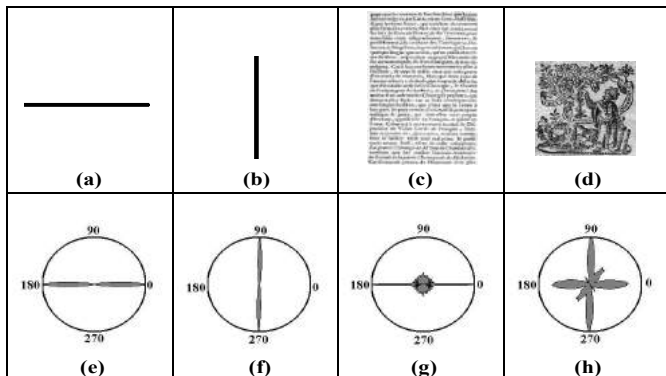
$$R_{(x,y)}^{I(\alpha,\beta)} = \sum_{\alpha \in \Omega} \sum_{\beta \in \Omega} I(x,y) I(x + \alpha, y + \beta) \quad (1)$$

$$= FFT^{-1}([FFT [I(x,y)] FFT^* [I(x,y)]]) \quad (2)$$

$$R_{(x,y)}^I(\Theta_i) = \sum_{D_i} R_{(x,y)}^{I(\alpha,\beta)} \quad (3)$$

$$R'_{(x,y)}^I(\Theta_i) = \frac{R_{(x,y)}^I(\Theta_i) - R_{min}^I}{R_{max}^I - R_{min}^I} \quad (4)$$

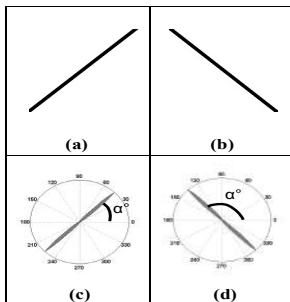
Examples of directional roses



{**(a),(b),(c),(d)**} are the original images,
 {**(e),(f),(g),(h)**} are respectively their roses of directions.

First texture feature : Main angle of the directional rose

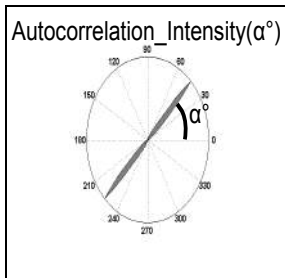
$$F_{(x,y)}^{(1)} = |180 - \operatorname{argmax}_{\Theta_i \in [0,180]} (R'_{(x,y)}(\Theta_i))| \quad (5)$$



⇒ Principal orientation information.

Second texture feature : Intensity of the autocorrelation function

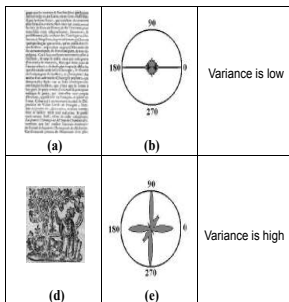
$$F_{(x,y)}^{(2)} = R_{(x,y)}^I(\operatorname{argmax}_{\Theta_i \in [0,180]} (R_{(x,y)}^I(\Theta_i))) \quad (6)$$



⇒ Level of anisotropy of the analysis window.

Third texture feature : Variance of the rose intensities

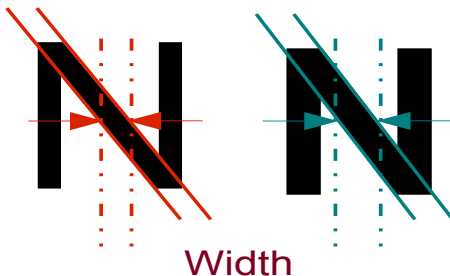
$$F_{(x,y)}^{(3)} = \sigma^2(R'_{(x,y)}(\Theta_i)) \quad (7)$$



⇒ Overall shape of the rose.

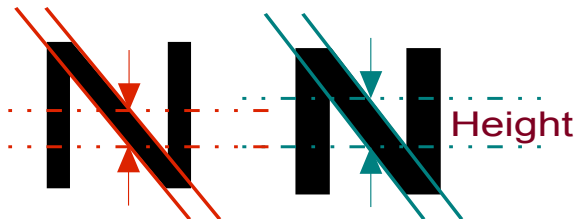
Fourth texture feature : Mean stroke width

$$F_{(x,y)}^{(4)} = \sum_{\Theta \in [10,80]} |I(x,y) - T_{(\alpha,0)}^{\Theta}(I(\frac{y}{|\tan(\Theta)|}, y))| \quad (8)$$



Fifth texture feature : Mean stroke height

$$F_{(x,y)}^{(5)} = \sum_{\Theta \in [10,80]} |I(x,y) - T_{(0,\beta)}^{\Theta}(I(x, x * |\tan(\Theta)|))| \quad (9)$$



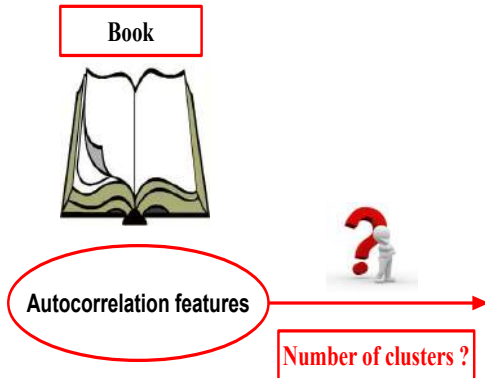
Layout segmentation and characterization of ancient document images

Book

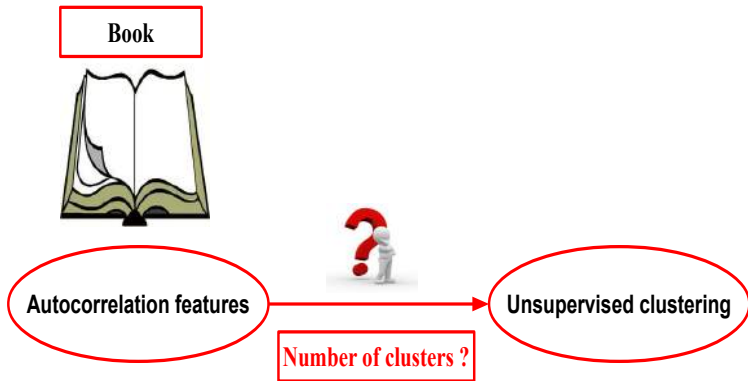


Autocorrelation features

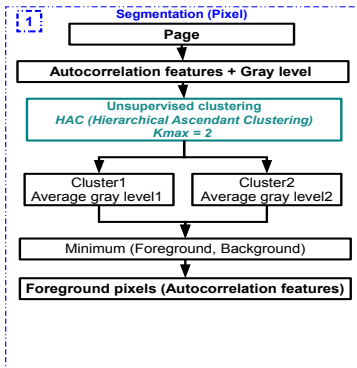
Layout segmentation and characterization of ancient document images



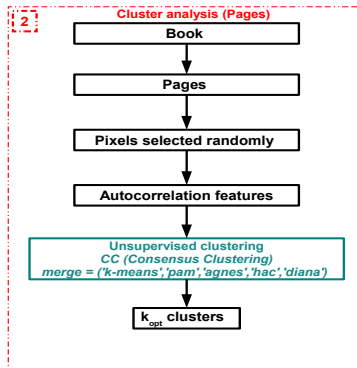
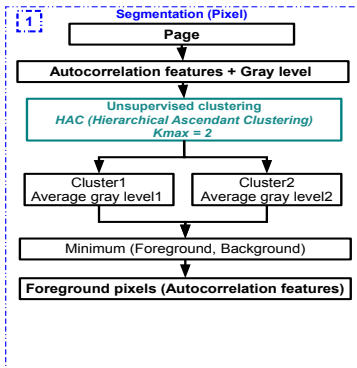
Layout segmentation and characterization of ancient document images



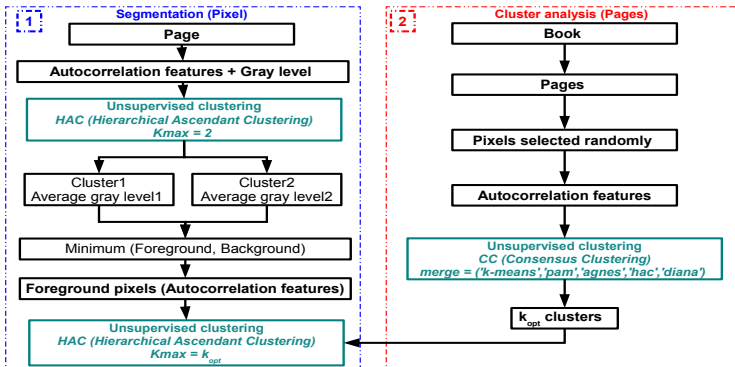
Layout segmentation and characterization of ancient document images



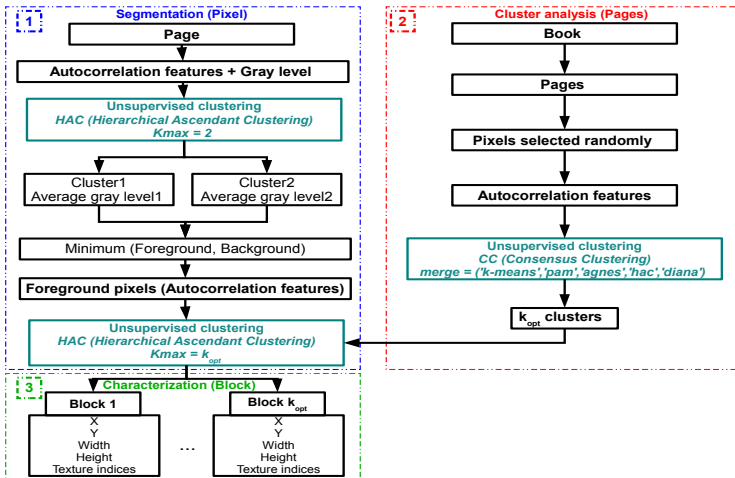
Layout segmentation and characterization of ancient document images



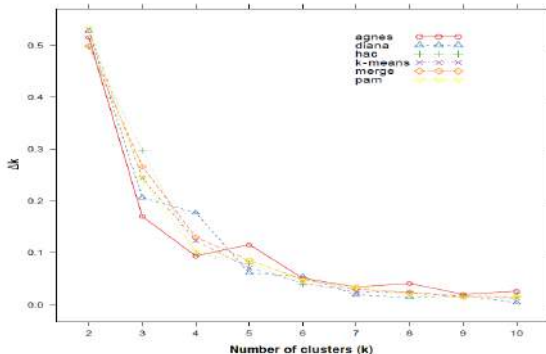
Layout segmentation and characterization of ancient document images



Layout segmentation and characterization of ancient document images



Consensus clustering



Plot of Δk changes in area under the cumulative density curve for the consensus matrix for each clustering experiment against cluster number k [11].

[11] T. Simpson, J. Armstrong, and A. Jarman, "Merged consensus clustering to assess and improve class discovery with microarray data," in *Boston Medical Center Bioinformatics (BMC Bioinformatics 2010)*. BioMed Central, 2010, pp. 1471–1482.

Outline

1 Method

- Previous work
- Proposed method
 - Segmentation using the autocorrelation function and multiresolution
 - Unsupervised clustering approach

2 Evaluation

- Corpus
 - Overview
 - Examples
- Experimental protocol
 - One font and graphics
 - Two fonts and graphics
 - Only two fonts
- Evaluation
 - Homogeneity measure
 - Example of segmentation and evaluation result
- Results
 - Examples of segmentation result
 - Accuracy metrics

3 Conclusion and further work

- Conclusion and further work
- Contribution

- 91596 images of ancient documents from 515 different books

- 91596 images of ancient documents from 515 different books
- Encompass six centuries (1200-1900) of France history

- 91596 images of ancient documents from 515 different books
- Encompass six centuries (1200-1900) of France history
- Grayscale and color images digitized with 300 and 400 dpi

- 91596 images of ancient documents from 515 different books
- Encompass six centuries (1200-1900) of France history
- Grayscale and color images digitized with 300 and 400 dpi
- 242 Grayscale books and 283 color books

- 91596 images of ancient documents from 515 different books
- Encompass six centuries (1200-1900) of France history
- Grayscale and color images digitized with 300 and 400 dpi
- 242 Grayscale books and 283 color books
- Images saved in tif format (3.4T)

- 152 printed monographs (1400-1900)

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)
 - 18 grayscale books and 128 color books

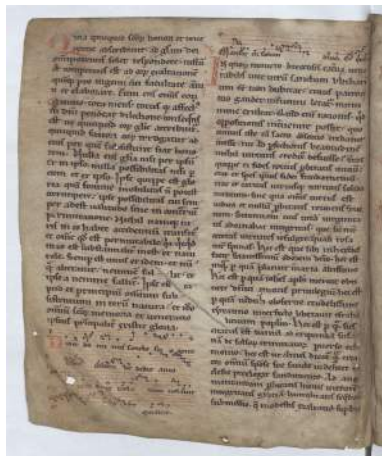
- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)
 - 18 grayscale books and 128 color books
 - 54406 images of documents

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)
 - 18 grayscale books and 128 color books
 - 54406 images of documents
- 217 journals and newspapers (1800-1900)

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)
 - 18 grayscale books and 128 color books
 - 54406 images of documents
- 217 journals and newspapers (1800-1900)
 - 194 grayscale books and 33 color books

- 152 printed monographs (1400-1900)
 - 30 grayscale books and 122 color books
 - 34532 images of documents
- 146 manuscripts (1200-1800)
 - 18 grayscale books and 128 color books
 - 54406 images of documents
- 217 journals and newspapers (1800-1900)
 - 194 grayscale books and 33 color books
 - 2658 images of documents

Gallica † : Manuscripts

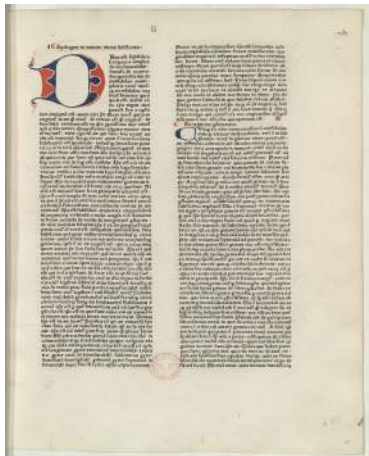


Experimental protocol

Gallica † : One font and graphics



Manuscript



Printed

Gallica † : Two fonts and graphics



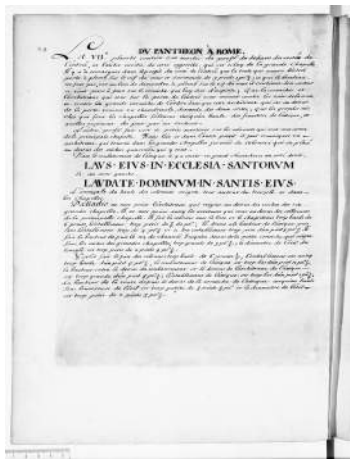
Manuscript



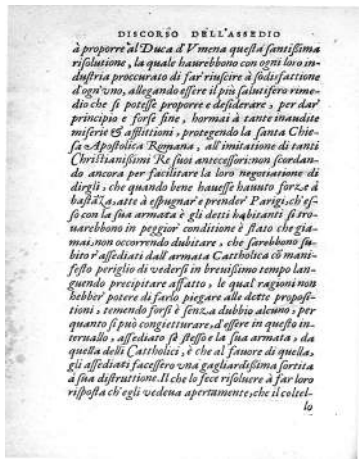
Printed

Experimental protocol

Gallica † : Only two fonts



Manuscript



Printed

Homogeneity measure

$$H(B, G) = \frac{1}{|G|} \sum_j \frac{\max_{1 \leq k \leq k_{opt}} (|b_i, (b_i \in g_j) \wedge (l_{B_i} = k)|)}{|\{b_i \in g_j\}|} \quad (10)$$

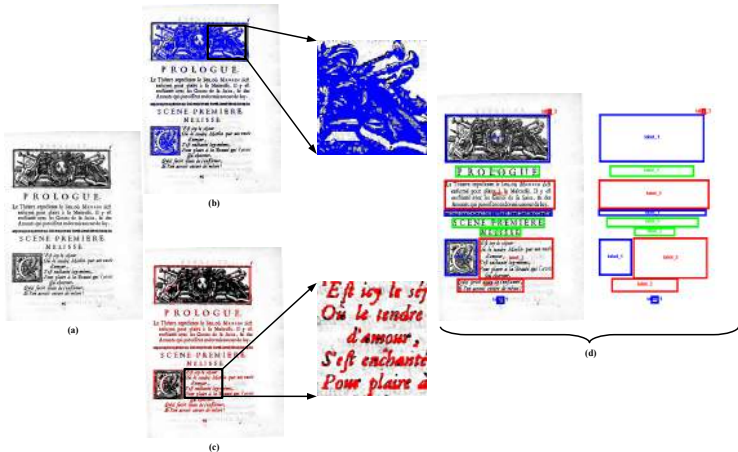
where $|\cdot|$ is the number of pixels in the given block.

$B = \{b_1, b_2, \dots, b_i, \dots, b_n\}$ and $G = \{g_1, g_2, \dots, g_j, \dots, g_m\}$ are respectively the sets of result blocks and rectangular regions of the ground-truth.

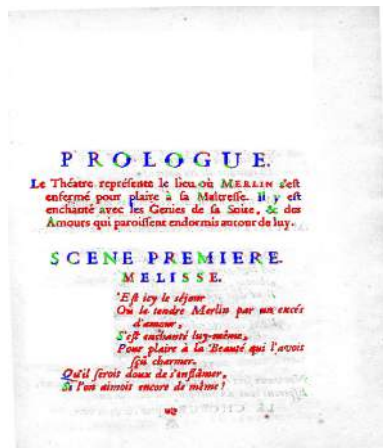
$L_B = \{l_{B_1}, l_{B_2}, \dots, l_{B_i}, \dots, l_{B_n}\}$ corresponds to a set of labels obtained with our clustering methodology.

Example of segmentation and evaluation result

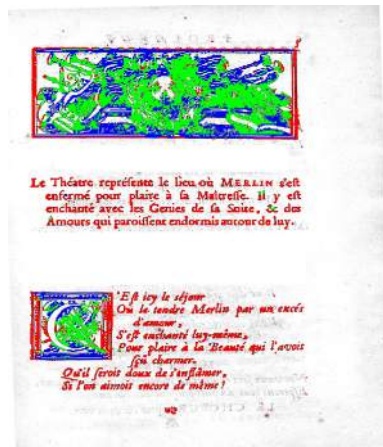
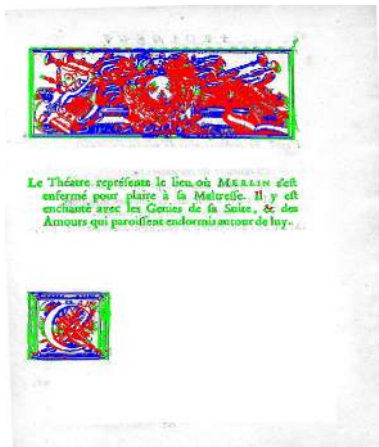
(a) original grayscale image, (b) cluster representing the graphics, (c) cluster representing the text and (d) ground-truth (Ground-truthing Editor (GEDI)).



Simplified images



Simplified images



Simplified images



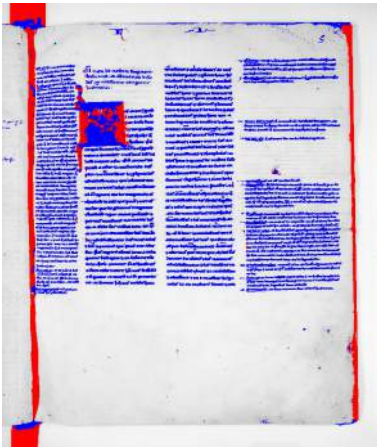
Printed-Two fonts and graphics



Manuscript-One font and graphics



Manuscript- Two fonts and graphics



Homogeneity measure

- H is based on spatial overlaps of the ground-truth rectangle and the segmentation result

	Document category	Document content	Pages	$\mu(H)$	$\sigma(H)$	
$H(B, G)$	Manuscript	One font and graphics	50	0,94	0,03	
		Two fonts and graphics	56	0,84	0,05	
		Only two fonts	50	0,87	0,05	
		Overall	156	0,88	0,04	
	Printed	One font and graphics	50	0,84	0,14	
		Two fonts and graphics	50	0,80	0,05	
		Only two fonts	60	0,80	0,10	
		Overall	160	0,81	0,09	
	Overall			316	0,85	0,07

$\mu(H)$, $\sigma(H)$, $Max(H)$ and $Min(H)$ are respectively the mean value, standard deviation, maximum and minimum values of the

Outline

1 Method

- Previous work
- Proposed method
 - Segmentation using the autocorrelation function and multiresolution
 - Unsupervised clustering approach

2 Evaluation

- Corpus
 - Overview
 - Examples
- Experimental protocol
 - One font and graphics
 - Two fonts and graphics
 - Only two fonts
- Evaluation
 - Homogeneity measure
 - Example of segmentation and evaluation result
- Results
 - Examples of segmentation result
 - Accuracy metrics

3 Conclusion and further work

- Conclusion and further work
- Contribution

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)
- Computation of statistical attributes (Grey Level Co-occurrence Matrix, Entropy, Homogeneity degree, Connection degree, etc.)

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)
- Computation of statistical attributes (Grey Level Co-occurrence Matrix, Entropy, Homogeneity degree, Connection degree, etc.)
- Computation of other texture features (Wavelets, etc.), and LBP, etc.

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)
- Computation of statistical attributes (Grey Level Co-occurrence Matrix, Entropy, Homogeneity degree, Connection degree, etc.)
- Computation of other texture features (Wavelets, etc.), and LBP, etc.
- Feature selection (Sequential Forward-Backward Search (SFBS), Genetic Algorithm (GA), etc.)

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)
- Computation of statistical attributes (Grey Level Co-occurrence Matrix, Entropy, Homogeneity degree, Connection degree, etc.)
- Computation of other texture features (Wavelets, etc.), and LBP, etc.
- Feature selection (Sequential Forward-Backward Search (SFBS), Genetic Algorithm (GA), etc.)
- Recursive clustering

Conclusion and further work

- Automatic, non-parametric, and unsupervised method for the segmentation and characterization of historical document images without any *a priori* knowledge
- Computation of frequency descriptors (Gabor filters)
- Computation of statistical attributes (Grey Level Co-occurrence Matrix, Entropy, Homogeneity degree, Connection degree, etc.)
- Computation of other texture features (Wavelets, etc.), and LBP, etc.
- Feature selection (Sequential Forward-Backward Search (SFBS), Genetic Algorithm (GA), etc.)
- Recursive clustering
- Definition of one or more signatures for each page

Contribution

- " Old document image segmentation using the autocorrelation function and multiresolution analysis ", M. Mehri, P. Gomez-Krämer, P. Héroux and R. Mullot, Document Recognition and Retrieval (DRR-XX), Part of the IS&T/SPIE 25th Annual Symposium on Electronic Imaging, 2013 February 5-7, San Francisco, CA, USA

Contribution

- " Old document image segmentation using the autocorrelation function and multiresolution analysis ", M. Mehri, P. Gomez-Krämer, P. Héroux and R. Mullot, Document Recognition and Retrieval (DRR-XX), Part of the IS&T/SPIE 25th Annual Symposium on Electronic Imaging, 2013 February 5-7, San Francisco, CA, USA
- " Extraction of metadata related to "image" and "structure" of old documents ", M. Mehri, P. Gomez-Krämer, P. Héroux and R. Mullot, INRIA-Visual Recognition and Machine Learning Summer School (VRML), 2012 July 10, Grenoble, France

References (1)

- [1] S. Bres, "Contributions à la quantification des critères de transparence et d'anisotropie par une approche globale : application au contrôle de qualité de matériaux composites," Ph.D. dissertation, Institut National des Sciences Appliquées de Lyon, 1994.
- [2] R. Ferrell, S. Gleason, and K. Tobin, "Application of fractal encoding techniques for image segmentation," in *International Conference on Quality Control by Artificial Vision VI (QCAV 2003)*. SPIE, 2003, pp. 69–77.
- [3] R. Haralick, K. Shanmugam, and I. Dinstein, "Textural features for image classification," in *Systems Man and Cybernetics (SMC)*. IEEE, 1973, pp. 610–621.
- [4] A. Jain and S. Bhattacharjee, "Text segmentation using gabor filters for automatic document processing," in *Machine Vision and Applications*. Springer, 1992, pp. 169–184.
- [5] N. Journet, J. Ramel, R. Mullot, and V. Eglin, "Document image characterization using a multiresolution analysis of the texture : application to old documents," in *International Journal of Document Analysis and Recognition (IJ DAR 2008)*. Springer-Verlag, 2008, pp. 9–18.
- [6] S. Khedekar, V. Ramanaprasad, S. Setlur, and V. Govindaraju, "Text - image separation in devanagari documents," in *International Conference on Document Analysis and Recognition (ICDAR 2003)*. IEEE, 2003, pp. 1265–1269.
- [7] S. Nicolas, Y. Kessentini, T. Paquet, and L. Heutte, "Handwritten document segmentation using hidden markov random fields," in *International Conference on Document Analysis and Recognition (ICDAR 2005)*. IEEE, 2005, pp. 212–216.
- [8] A. Ouji, Y. Leydier, and F. LeBourgeois, "Chromatic / achromatic separation in noisy document images," in *International Conference on Document Analysis and Recognition (ICDAR 2011)*. IEEE, 2011, pp. 167–171.

References (2)

- [9] T. Pavlidis and J. Zhou, "Page segmentation and classification," in *Graphical Model and Image Processing (CVGIP 1992)*. Elsevier Science, 1992, pp. 484–496.
- [10] C. Sabharwal and S. Subramanya, "Indexing image databases using wavelet and discrete fourier transform," in *Symposium on Applied Computing (SAC 2001)*. ACM, 2001, pp. 434–439.
- [11] T. Simpson, J. Armstrong, and A. Jarman, "Merged consensus clustering to assess and improve class discovery with microarray data," in *Boston Medical Center Bioinformatics (BMC Bioinformatics 2010)*. BioMed Central, 2010, pp. 1471–1482.
- [12] M. Tuceryan, "Moment based texture segmentation," in *Pattern Recognition Letters*. Elsevier Science, 1994, pp. 659–668.
- [13] M. Tuceryan and A. K. Jain, "Texture segmentation using voronoi polygons," in *Pattern Analysis and Machine Intelligence*. IEEE, 1990, pp. 211–216.
- [14] K. Wong, R. Casey, and F. Wahl, "Document analysis system," in *IBM Journal of Research and Development*. IBM Research Division, 1982, pp. 647–656.

Thank you for your attention

*** Thank you for your attention ***

Thank you for your attention

*** Thank you for your attention ***

Thank you for your attention

*** Thank you for your attention ***

Thank you for your attention

*** Thank you for your attention ***