

WP3-DIGIDOC

Thesis: Extraction of metadata related to "image" and "structure" contained in old documents

Prepared by: Maroua MEHRI^{1,2}

Supervised by: Petra GOMEZ-KRÄMER¹

Pierre HEROUX²

Rémy MULLOT¹

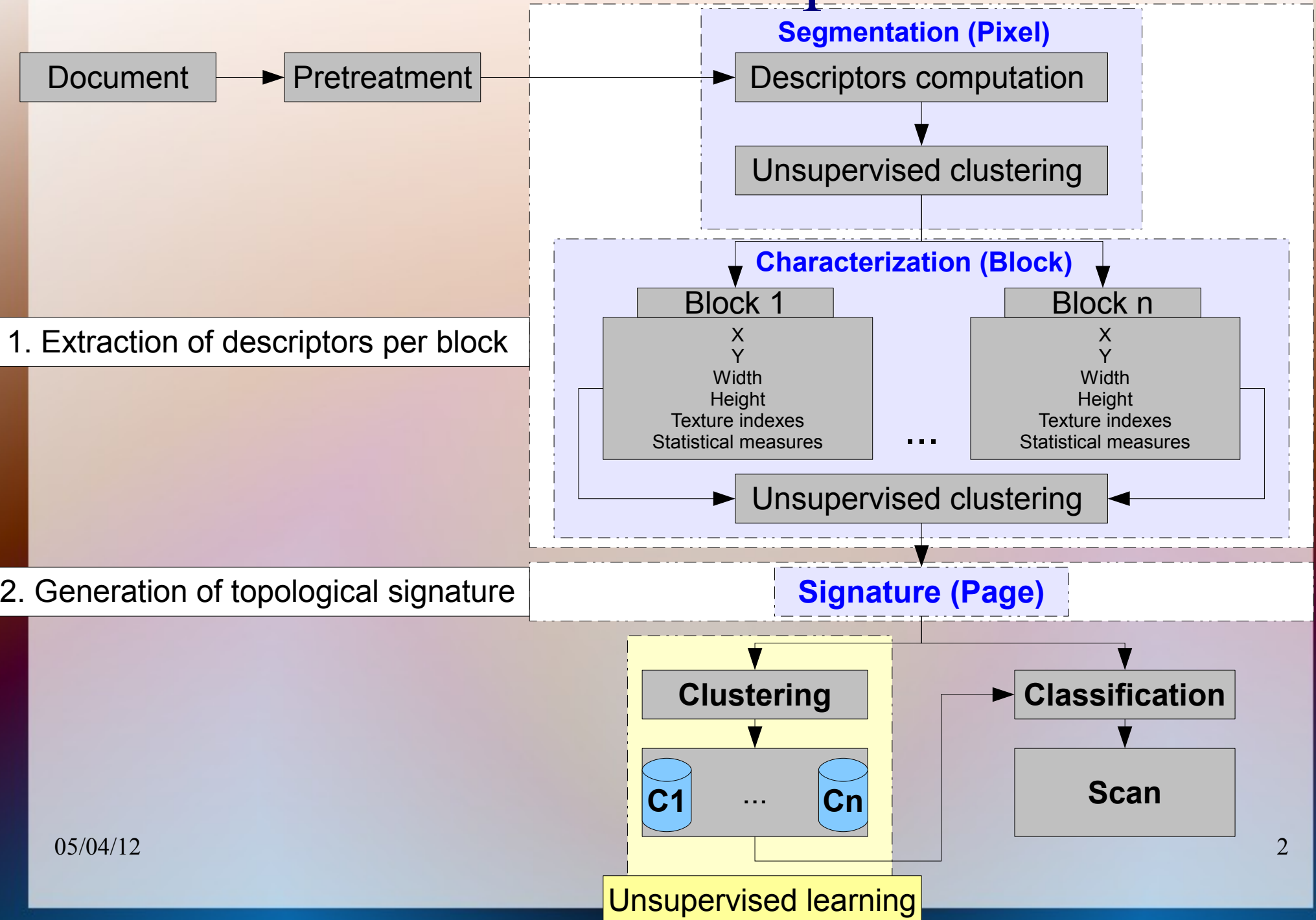
[1] L3I, University of La Rochelle, Avenue Michel Crépeau, 17042 La Rochelle, France

E-mails: maroua.mehri@univ-lr.fr, petra.gomez@univ-lr.fr, remy.mullot@univ-lr.fr

[2] LITIS, University of Rouen, Avenue de l'Université, 76800 Saint-Etienne-du-Rouvray, France

E-mail: pierre.heroux@univ-rouen.fr

Research topic



Texture attributes of autocorrelation

- **Angle: A_{ij}**
 - Corresponds to the main angle of the rose of directions
- **Autocorrelation intensity: I_{ij}**
 - Corresponds to the intensity of the autocorrelation function for the main orientation found (1st index)
- **Standard deviation: Std_{ij}**
 - Corresponds to the standard deviation of the intensities of the rose, expect for the orientation of maximal intensity
 - $Std_{ij} \ll 1$: The main orientation is significantly more prevalent than other orientations
 - $Std_{ij} \gg 1$: The rose of directions is deformed and a large number of orientations are present in different proportions [Jou06, JRMV08]

Texture attributes of autocorrelation

$$\begin{aligned}\mathcal{D}^h(I)_0 &= 0 \\ \mathcal{D}^h(I)_n &= \sum_{(x,y) \in \Omega} \|I(x,y) - T_h(I,n)(x,y)\|\end{aligned}$$

With $T_h(I, \delta)$ is the translation of the gray-scale image I (defined on the plane Ω) of δ pixels along the horizontal axis.

- **Mean stroke width: S_w**

- Corresponds to an estimation of the mean stroke width of an image computed from the sequence $\mathcal{D}^h(I)_n$ until $n = n_m$, when its growth rate becomes lower than 10%. Then, we obtain $S_w = n_m$

- **Mean Stroke height: S_h**

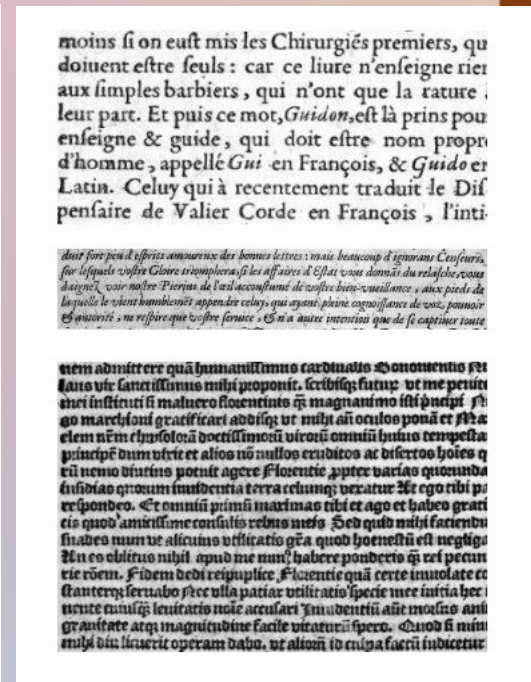
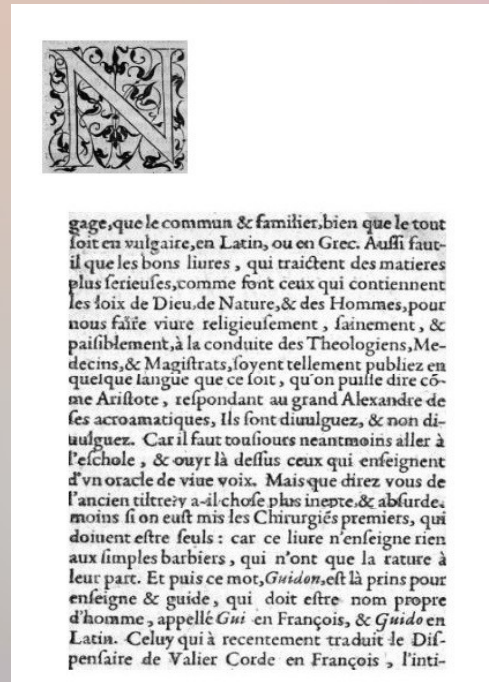
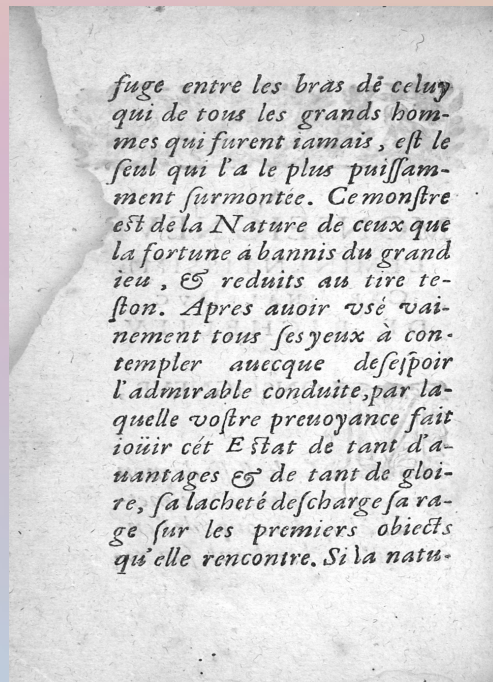
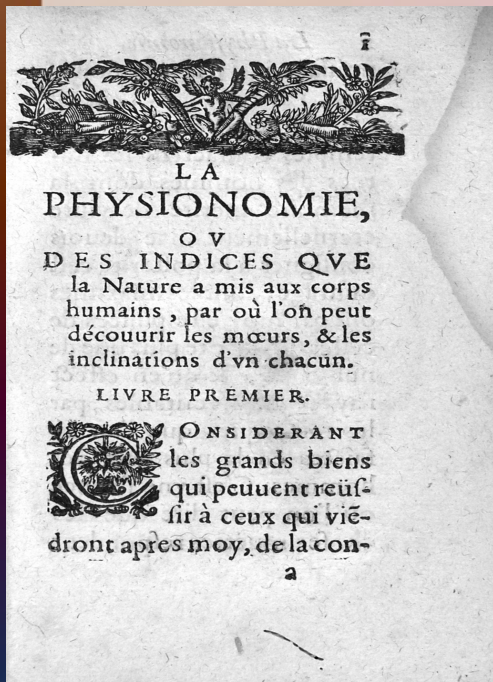
- Corresponds to an estimation of the mean stroke height of an image computed from the sequence $\mathcal{D}^v(I)_n$ until $n = n_m$, when its growth rate becomes lower than 10%. Then, we obtain $S_h = n_m$ [OLL11, OLL12]

Texture attributes computation

- Assigning a feature vector for each given number of pixels of the image using a sliding window
- Fixed-size window, positioned on the left corner of the image
- 5 texture attributes of autocorrelation are calculated
- The point at the center of the window is then marked by the 5 numerical values calculated
- Window is moving from a number of pixels and so on until going through the whole image

Database

- 25 images of old documents
 - 16 images containing text and graphic areas (drop cap letters, ornaments...) including 7 images are snapshots
 - 8 images with only text zones (text containing many paragraphs characterized by different fonts size) including 7 images are snapshots
 - 1 image with only graphic block



Examples of database images

Constraints

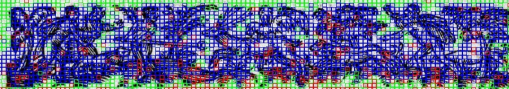
- Effectiveness (segmentation and clustering performance)
- Efficiency (computing time)
- Number of clusters
- Large dataset for clustering
- Problem of out of memory

Parameters

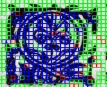
- **Multi-resolution**
 - Keep the size of the original image and varying the size of the analysis window
- **Size of window**
 - [8, 16, 32, 64]
- **Pixels number of window shift**
 - 8
- **Clustering method**
 - *HAC (Hierarchical Ascendant Clustering)*
 - *Consensus clustering*

Results

32

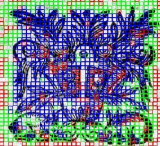


L. IOUBERT,
AV LECTEUR



Q V A N D j'ay promis des annotations fort amples sur la grande Chirurgie de M. G. V. j'entendois qu'on imprimait le Latin premierement, & que le François viendroit apres. Mais le Libraire en est allé au delà d'une copie, & en a d'autre aduis, & a voulu commencer par le François: dont les dites annotations n'ont esté si facilement traductibles, que l'on ne s'est esté de l'imprimer. Aussi de l'un de ceux en ça j'ay esté fort delivré de cette besongne, pour avoir vacqué longuement au service du Roy, & du Roy de Navarre: Mais j'espère dans peu de mois l'achever entièrement à ma promesse. Cependant on jouira de cette Chirurgie, mieux traduite que n'a esté parcy devant: & aussi tost que la Latine, par moy corrigée (qui est maintenant sous la presse) avec mes annotations en même langue, auront vu la lumière, les dites annotations en François se trouveront prestes à imprimer. Dieu aydant, auquel seul en soit la gloire, & le profit à tout le monde.

CV



LES RECEPTES.

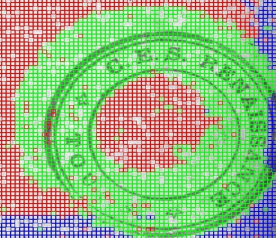
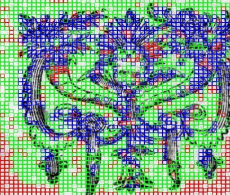
Syrop de Guy d'igerant le phlegme. 629. 18

T

T Abille à rinder l'umbr. 626. 16
Tetrastarmac. 627. 13
Theriacque de singes fect. 547. 29
Triastarmac de Calen. 203. 29

Trochiscs pour la douleur de dents avec chaleur. 547. 20
Trochiscs pour la goutteuse. 592. 17
Trochiscs pour la furdie & im-
riement. 552. 10
Trochiscs aldar. 580. 36
Trochiscs d'apoplexies. 580. 30
Trochiscs caluicem. 681. 5
Trochiscs de berberis. 593. 15
Trochiscs de chalc. 680. 26
Trochiscs de lacque. 195. 33
Trochiscs narcomen. 668. 26

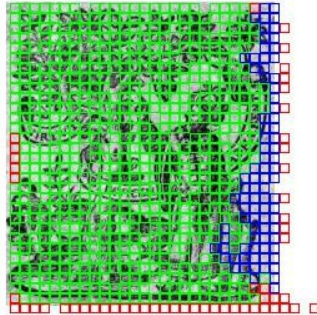
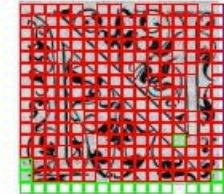
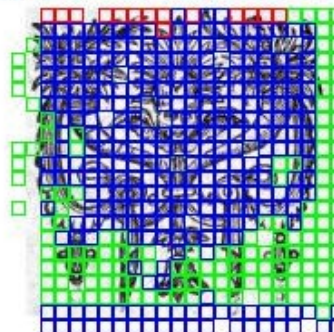
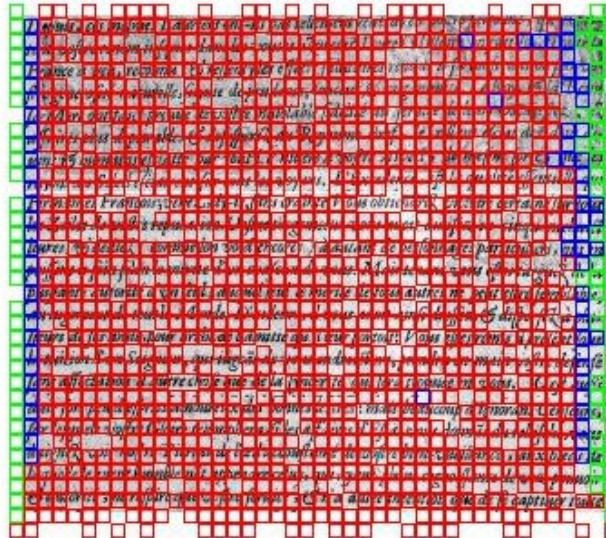
F I N



Privilege du Roy.

L OVIS par la grace de Dieu Roy de France & de Navarre. A nos Amex & feaux Con-
seillers. Les gens tenans
nos Cours de Parlement,
Maistres des Requestes Or-
dinaires de nostre Hostel,
Preuost de Paris, son Lieute-
nant, & autres de nos Juges
qu'il appartiendra, Salut.
Nostre bien aimé Henry de
Boyun du Vauroit âgé de
XII. ans, nous a tres hum-
blement fait remonstrer

Results

[illegible]

gare, que le commun & familier bien que le tout
soient en vulgaire, en Latin, ou en Grec. Aussi faut
il que les bons iurés, qui traitent des matieres
plus erudites, comme font ceux qui contiennent
les loix de Dieu de Nature, & des Hommes, pour
nous faire vivre religieusement, sainement &
pudiquement, à la conduite des Theologiens, Me-
decins, & Magistrs, ne soient tellement publiez en
quelque langue que ce soit, qu'on puisse dire, se-
lon Aristote, respondant au grand Alexandre de
ses aecademiques, Ils l'ont en vulgaire, & non en
vulgaire. Car si l'on confond neantmoins avec
l'echole, & ouyrla plus de ceux qui enseignent
en grande de vaine voix. Mais que dire, vous di-
l'ancien vulgaire a l'histoire plus inepte & absurde
moins s'en eust mis les Chirurgiens premiers, qui
doient estre seuls, car ce livre n'enseigne rien
aux simples barbers, qui n'ont que la rature a
leur part. Et puis ce mot, *Quoniam*, est la pris pour
enleigne & guide, qui doit estre non propre
à l'homme, appellé *Gai* en François, & *Quoniam*
Latin. Celui qui a recentemente traduit de l'Es-
paignois de Valer *Quoniam* en François, l'anti-

Future works

- **Consensus clustering and validation of autocorrelation descriptors**
- **Frequency attributes**
 - *First feature*: Characterizes the transitions between paper and ink
 - *Second feature*: Characterizes the white spaces separating the collateral elements
- **Statistical attributes**
 - *Grey Level Co-occurrence Matrix*
 - *Entropy*
 - *Homogeneity degree*
 - *Connection degree*
- **Texture attributes**
 - *Wavelets*

References

- [Jou06] Journet, N. Analyse d'images de documents anciens : une approche texture. PhD thesis, L3i, 2006.
- [JRMV08] Document Image Characterization Using a Multiresolution Analysis of the Texture: Application to Old Documents. N. Journet, J.Y Ramel, R. Mullot, V. Eglin. International Journal of Document Analysis and Recognition (IJDAR 2008).
- [OLL11] Chromatic / achromatic separation in noisy document images. A. Ouji, Y. Leydier, F Lebourgeois. International Conference on Document Analysis and Recognition (ICDAR 2011), Beijing, China. 2011.
- [OLL12] Extraction de texte à base de segmentation colorimétrique dans les images de presse. A. Ouji, Y. Leydier, F Lebourgeois. CIFED, Bordeaux, France. 2012.

Thank you for your attention