



# From pixels to content: An overview of the main techniques used in DIA

#### **Jean-Yves RAMEL**







#### From Pixels to Contents Introduction





#### From Pixels ... to contents Outline

- From pixels...
  - □ What is an image?
  - □ Image (pre-)processing

#### ... to Text

- □ Transcription and Layout analysis
- Segmentation and content extraction
- An overview of Pattern Recognition
- ... but also to non-Text
  - Content characterization and signatures
  - Content retrieval and spotting
- Back to meta-data?
  - From descriptive to perceptual meta-data
  - □ Is there adequate encoding formats?
- Conclusions and perspectives

# From Pixels...





#### 

- Images come from a grid of microscopic photosensitive cells called **PIXELS**
- Sampling



#### Quantization

- □ Assignment of a numerical value drawn from the received lighting energy / pixel (grid unit)
- □ Continuous value  $(xi,yi) \rightarrow Discrete value (xi,yi) \rightarrow Pixels$
- The range of colors that each pixel can take









#### From Pixels... What is an image?

#### Image Quantization

**Binary images:** l(i,j) = 0 black or l(i,j) = 1 white

Gray level (8 bits/pixel) images: I(i,j) = 0....255 from the lighter to the darker.

Color images (24 bits/pixel): 3 values of lighting intensity Red, Green, Blue

 $I_1(i,j) = 0 \dots 255 - I_2(i,j) = 0 \dots 255 - I_3(i,j) = 0 \dots 255$ 

Image Representation & Processing

Image = Array(s) of pixels = Matrices of values 1 pixel = A position inside the image (i,j) + 1 color (1 to 3 values)

The values *l(i,j)* associated to each pixel *s(i,j)* represent their brightness intensity



32 bits

1 bit



~	2			Niveaux de gris - 8 bits: 0 - noir 255 - blanc					
6	64	60	69	100	149	151	176	182	179
	65	62	68	97	145	148	175	183	181
ALL ALL	65	66	70	95	142	146	176	185	184
0.000	66	66	68	90	135	140	172	184	184
210	100	64	64	84	129	134	168	181	182
	59		62	88	130	128	166	185	180
	60		60	85	127	125	163	183	178
	62	62	58	81	122	120	160	181	176
	63	64	58	78	118	117	159	180	176



#### From Pixels... What is an image?

#### Sampling → Image resolution

- Number of pixels per lenght unit
- In dpi (dots per inches) or ppp (points par pouce)
- When the resolution decrease, the precision decrease





#### • VF Image processing

- Page A4 = 21x29.7 cm
- 200 dpi : 1650 x 2340 pixels = 3 861 000 pixels
- 300 dpi : 3500 x 2480 pixels = 8 680 000 pixels
- 16M colors, 1 pixel = 3 octets → 10 à 25 Mo/page !
- A trade-off between quality-quantity/time is mandatory
- Fidelity of the numerical version
- Mass of storage size Transmission / Processing time



#### From Pixels... Why few pixels are so important?



lated patterns can correspond pixels) !

on the boundaries of the shapes ocess

arance of characters until the

acters or touching characters **e resolution** 



#### procedure and ments of patterr that in spite of

OCR errors



#### From Pixels... What is Image (pre-)processing?

After the digitization, the images usually still have a lot of defaults

- Curvature and skew due to scanning
- Noise on boundaries, dots, blur, ...

9





# From Pixels... Image (pre-)processing

#### Curvature and skew correction is possible on text images

point & le refte de fàtin blanc, & tout paffemente & pou filé d'or : celluy des Efpingliers bonnet, collet, chautles, fouliers de uelours noir le pourpoint de fatin cramoify, doubleure des chauffes correspondât, rayez de paffemen & traifes d'or. Apres lefquelz paffoient quelques premie rangs armez & accompaignez de deux centz & fept Tif rans portantz rouge & noir les troys Enfeignes derrie euls braues & bien en ordre, & marchantz deu at deux cer cinquâte fix Cordoanniers ueffus de blanc & noir, laiffar à leurs et paules les troys Lieutenantz autant brauement ordre, & conduifantz centz quatre uingtz & douze Efp gliers portantz le pourpoint de uelours, fatin, ou taffe rouge, le collet & bonnet noir auec plume blanche, & zgri fariffaifant à chafcun,

Tout d'un ordre furuint la fixiefine Bande autant bel que plaifante pour la diuerfité des couleurs: la quelle coms ca par le rang de les troys Capitaines de Rue neuve acce firé de uelours noir, blanc, & bleu mouchetté menuemi de bouttons d'or, accópaigné du Capitaine des Chappeli ueitu de uelours blanc & noir & uerd à petitz grains d' fuyuant d'on melme pas aueccelluy des Fondeurs en ha de uelours blanc, & noir, & aurangé, recamé & bifetté d gent. Et lequel rang auec fes Tabourins & Fiffres de mél fur iuvuy d'aucurs autres armez de corfelezz & animes; & fuytre de Rueneuue en liuree de noir blanc & bleu; & nombre de quatre centz uingt & troys: lefquelz effoient fiez de troys Enfeignes fuyuantz auec melmes couleur leurs enfeignes, guidantz apres enha cent foisante & d point & le refte de fatin blanc, & tout palfementé & pou file d'or : celluy des Efpingliers bonnet, collet, chaudes, fouliers de uclours noir:le pourpoint de fatin cramoify, doubleure des chautles correspondât, rayez de passemen & trailes d'or. Apres lesquelz passoint quelques premie rangs armez & accompaignez de deux centz & lept Tif rans portantz rouge & noir:les troys Enleignes derrie culs braues & bienen ordre, & marchantz deuát deux cer cinquâte fix Cordoanniers uestus de blanc & noir, lassfar à leurs espaules les troys Lieutenantz autant brauement ordre, & conduisantz centz quatre uingtz & douze Esp gliers portantz le pourpoint de uelours, latin, ou tasse rouge, le collet & bonnet noir auec plume blanche, & gr: fatisfaifant à chascun.

Tout d'un ordre furuint la fixiefine Bande autant bel que plaifante pour la diuerfité des couleurs:laquelle cóm ca par le rang de les troys Capitaines de Rue neuue acce itré de uelours noir, blanz, & bleu mouchetté menuems de bouttons d'or, accopaigné du Capitaine des Chappeli ueltu de uelours blanc & noir & uerd à petitz grains d' fuvuant d'un mefine pas auec celluy des Fondeurs en ha de uelours blanc, & noir, & aurangé, recamé & bifetté d gent. Et lequel rang auec fes Tabourins & Fiffres de mel fut iuvuy d'aucuns autres armez de corfeletz & animes, & fuytte de Rueneuue en liuree de noir blanc & bleu, & nombre de quatre centz uingt & troys: lefquelz effoient ftez de troys Enfeignes fuyuantz auec mefines couleur leurs enfeignes, guidantz apres oulx cent foisante & f Chappellier de blanc noir & uert: Et à la fie les troys Lj



#### From Pixels... Image (pre-)processing

The problem is more complicated in case of heterogeneous content





#### From Pixels... Image (pre-)processing

The problem is more complicated in case of heterogeneous content



# ... to text and layout





Dans la Mer d'Inde il y a une espèce de poissons qui ont dans leur peau des poil longs que les gens, quant il les ont attrapés en font des vêtements pour s'habiller.





# From pixels ... to text Automatic Trancription

Manuscrit	Translation	]
R lamer dyndr 01 a une manier dynillons qui ont en leur praus prus fi lons quoles gen; en 3 05 font uchteurs p euls 3	Dans la Mer d'Inde il y a une espèce de poissons qui ont dans leur peau des poils si longs que les gens, quant il les ont attrapés, en font des vêtements pour s'habiller.	So many possibilities for transcriptions !
Diplomatic Transcription	Modernize une al	baleste lehue a cause de la
EN lamer dynde a une maniere de poi∬ons qui ont en leur piaus peus ∫i lons que les genz en font ue∫teures p̃ euls ue∫tir qñt il les ont p <sup>i</sup> s	En la mer a une man de poisson en leur pia lons, que l font vestet vestir quant il les ont pris.	et grantite dessus estempte

#### From pixels ... to text Automatic transcription but also layout Analysis

Afrans appe pluf to read floor on fringed angot aroftim le fame as forthe 2.13 In bolon of 17/6-2.6 Congright & Ca Diofer Call

Sum Ind finnerst robur round a New Allaloge & Gue a range & Ca Durk forme Ind trayour amor Ca filmer at get the Seffel refresphe

Form de glande uptoffer al long and han de registerion gome ch mile at fine In rommer & Cada alla fifter nature for Gay Capet for and high and my art for a grand Affer lang and my art for a grand forme on your officers of Grand and Suffer by art of grand and soffer by art of grand and and and and art of grand and and art officers of grand and and art officers of grand and and and and and art of grand and and art officers of and and art officers and and a grand and art officers of and a grand art officers of a grand art officers of and a grand art officers of a grand art officers officers of a grand art officers of a grand art officers of a grand art officers officers of a grand art officers of a grand art officers of a grand art officers officers of a grand art officers offic

7-receu dudit franceys rotier pour
8-une albaleste dehue a cause de la
9-dite ferme dudit trezein avec la
10-somme et quantite dessus estempte ?
11- Il flor(ins)

IIIIxXX V

12-receu de Alexandre ? escoffier al loup ? 13- dudit lieu de chatellion pour la 14-rense et ferme du commun de la dite Preservation of the communaute ? pour ung an commettre ? Iink between the uatre cent et cinquante? A ly par le Exerct and the Image au plus offrant du voloyr et 20-et consentement de pluseures des bourgeys 21-de la dite ville de chastellion apres 22-pluseures rues ....refaytes es lieus ?



As a first step, the segmentation result will have high impact on the final results provided by the OCR !!!!

#### « it is always the segmentation that raises problems »

- The segmentation is an irreversible processing because it is the result of an analysis of the image according to a criterion and a method.
- Many segmentation methods have been proposed:
  - □ Top-down methods (image decomposition into smaller and smaller elements) → Layout first characters after
  - □ Bottom-up methods (from pixels tp higher and higher level elements) → Characters first – layout after
  - Etc...



FoC

#### From pixels ... to text An overview of OCR mechanisms

#### First step : Image segmentation

- Transformation of the image (set of pixels) into patterns (regions of interest) of higher level (EoC)
- These EoC could be very simple (part of characters) or more sophisticated ones (paragraphs, illustrations, ...)
- EoC extraction: Background (white) / Foreground (black) separation





#### Just to illustrate the difficulties...

Most of the segmentation methods need a binarisation





Sonnet for Lena

O dear Lona, your beauty is so vast

It is hard sometimes to describe it fast.

If only your portrait I could compress.

And for your lips, sensual and tartual Thirtrens Grays found not the proper fractal.

Alms! First when I tried to use VQ

I thought the entire world I would impress

I found that your checks belong to only you.

Hard to match with sums of discrete cosines.

I pught have fixed them with backs here or there

Your silky hair contains a thousand lines.

# Just to illustrate the difficulties...

19 19

- Most of the segmentation methods need a binarisation
- Global threshold ->

#### Local thresholds ->

Niblack :  $S = m + ks^2$  avec k= -0,2 | *m* : mean et s : standard deviation

# <section-header><section-header><section-header><text><text>

Sonnet by a



#### Sonnet for Lena

O dear Loos, your beauty is so vast it is hard sometimes to describe it fast. I thought the entire world I would impress If only your portrait I could compress. Alast First when I tried to use VQ I found that your checks belong to only you. Your silky hair contains a thousand fines Hard to match with sums of discrete cosines. And for your lips, second and tactual Thirteen Crays found not the proper fractal. And while these setbacks are all quite server I might have fixed them with hacks here or there but when filters tooks parkle from your cyse I said, 'Dann all this. TB just digitize.'

Theras Caltharat



#### First step: Image segmentation / Connected components

Then, we can try to group black pixels together to localize and recognize higher level Element of Content (EoC)

Q, Q 1, 5	Secundum flouiæ iatus.	
Yoy 1,3	Prima tibiæ oßis linea. d	Л;
E, E1, 2, 3,4	Tertiatibiæ oßis linea.	
#1,3	Secundum tibiæ oßis latus.	
n 2,4	Afbera tertij tibiæ oßis lateris linea, cui musculus	il



N TIBIA limiliter atque in cu crates perpetuò tibiæ offa appella quod totius membri nomine lui exterius locatur, & interiori craffi eous, Latinis autem fura & fibula i dio craffius os, tibiæ os appellabo guens, quod his offibus, musculi

Tibiæ oßis o fibula ap. pendices. Sedes cuife,

nia ambienti formatur. Gracilius autem os fibula opt ci, quia ueluti tibiæ uenter eft, yasponunulau nuncupan suprà infracta appendix coalescit, ac superior quide tibi eft, anterioritamen sede crassior, ac in anteriori tibiære fum duci cernitur. Huic\* duo oblongi insculpuntur l mur ad tibiam distincti, & lubrica cartilágine incrustati. His sinibus i ANNICHTAN







#### Next step : Layout analysis

- Connected Components → Words → Lines → Paragraphs → Page
- The results have to be saved in XML format (Alto, ...)
- Choosing how to organize the XML tree (physical / logical) is not so easy...

	Newsweek
	Initial Pedidi Prove Tinisi Med Bataka Reenveloring Hell In Leditor Vienneti U.S. Affairs The Radical Windon by Ine Klein A Powell Source by Missed Fit Why Powell's Base Mattern by Josephan Aler Antifice A Law of Their Own World Affairs
How to Capture America's Radical Middle Car formed radie the 'warship because the former select that Areasines's angle has been particularly and the target of most inportent former in the powerful and anguage of AV Midle works bok for how. No ware rate to be being the other senses that polarize many set the sense of the best of the sense that any even the sense of the powerful and the set of the best polarize many set the sense of the sense of the sense of the With Blaux.	Wenner, Alvison Up Trout by Ca Intred Fore Robot for Tox Lines Units: Don't Cry for Meson Renge: Wall Hard Kill the Third? Isla India: Fight for the Minach Tree by Mithand Toxic Pathone A Constant Ad Digits Officer A Constant Pathon Officer A Constant
Deliver Us From Evil Is Runai's Urals, gold re- mann. They defloated a former labor cares to the memory of the Secket system's visiting, in hopes that used an oxtrage value the secket system's visiting of the secket system's visiting of the secket system's visiting of the secket system's visiting of the secket system's visiting of	Technic U.S. Fock in Cyberrower Frages, Privos and Perdite by Robert J. Samuelson Scientify & Then Artis Techniculars: Theore Citics by In- Techniculars: Theore Citics by In- Techniculars: Theorem Citics by In-
Hiddand Wigner's four-port. The Biret rithe Nikolang 's Bihans of loval market is all controls of the Second Second right German Source comparison of the Work of the Thing's no reco- posed for work of the Thing's no reco- tained for work of the Thing's no re- tained for a second second second second second to be readed as the second second second second second to be readed as the second second second second second to be readed as the second second second second second to be readed second second second second second second to be readed second second second second second second second to be readed second second second second second second second to be readed second se	Epitode: Ceorge Previo E. Will B. Contra II. Reachance by the level Wills. The Contra II. Reachance by the level Wills. The Contra II. II. Sector Based by Reachance Based Based Based Based Based Based Based Based Based Based Based The Contra II. II. Sector Based Based Based Based Based Based Based The Contra II. Sector Based Based The Contra II. Sector Based Based Based Based Based Based Based Based Based The Contra II. Sector Based B







#### Next step : Layout analysis

Two kind of structures have been identified by researchers in DIA:

- The logical structure → the generic one corresponding to a priori knowledge about the content of the document
- The physical structure → the analysed instance corresponding to the extracted EoC inside the image, each one associated to descriptive features (size, position, number of sub-patterns, ... )
- Layout analysis tries to recognize these 2 structures (EoC identification)





#### Next step : Layout analysis

 The analysis / identification of the EoC is usually achieved based on a rule based system defined through a grammar (static one) or defined interactively by the users





Faille de la base d'apprentissage

Nombre de

Nombre de caractéristiques

Nombre de classes

100

50

2

100

caractéristiques

**Statistique suffisante** 

#### From pixels ... to text An overview of OCR mechanisms

#### Many possible choices and techniques

For selection of discriminative features Performance du classifieur 1 EoC  $\rightarrow$  1 Vecteur  $\vec{x} = \begin{vmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{vmatrix} \in IR^{-n}$ Malédiction de la dimensionalité Performance du classifieur 1 EoC → 1 Graph 360.62 15 Complexité vs nombre de caractéristiques 105.62 Complexité Machine Learning models and tools Nombre de classes Performance du classifieur

☆



#### Next step : Pattern / EoC recognition (toward Machine Learning)

#### How computers can recognize objects?

- We need a large set of (labelled) examples similar to the patterns to be recognized → a training set
- We need a list of stable and discriminative features (shape, color, size,...) used to describe the patterns (labelled ones and unknown one)





#### Next step : Pattern / EoC recognition (toward Machine Learning)

#### How computers can recognize objects?

 When an unknown EoC arrives, we compute its features and compare it with the content of the training set (associated built models)





#### Deep Learning (Conv. Neural Net)





#### Why commercial OCR are not working well on historical documents?

- Noises and degradations
- Unusual layout
- Unsuited training set



#### **Fine Reader**



# AA

Vena calla por tio à dextro cordis finu ad ingulu ufg conscendens, qua fi brarum in uenarii corpore natura do Arina gratia Studiole finximus. Si forte interim rds morum bic pasim abtruncatorum ra tione expetis, banc Affeititia & uende portione innon peculiaris tegra cana uena uena tuni, delineationi quinti

nis administrationi utiles fore cognouit. Q ad fe admi fanguinen ad fubleq fioneopu ftatim uni ram ex ali brasuena ribus uen nicæ conf appolitec rem tunic poris part

ca. Capitis fini subijcienda confer, in illa quid D.F.G.H.I.K.N.S.T. o a indi go tramite





# Why commercial OCR are not working well on historical documents?

- Lack of data, knowledge and experiences
  - $\hfill\square$  Unusual fonts and characters  $\clubsuit$  training data needs to be created
  - □ Unusual languages → Lexicons, dictionaries and language models need to be created
- Context often allows to modify our understanding of what is perceived by our senses
  - □ Until now, we tried to recognized EoC without using their context
  - □ The same EoC could be interpreted differently according to its surrounding context
  - □ Results of OCR are highly correlated to the adequacy of the used **word dictionary**



- Is there methods that need less a priori knowledge?
- Processing non-Textual parts can be good source of inspiration?

# ... to non-text





#### From pixels ... to non-textual contents Pictorial Content is also of high interest







Figures (+ de 1500)



#### BATYR: http://www.bvh.univ-tours.fr







#### Perceptual meta data instead of classical meta data

Computation of signatures for all the images or even sub-parts of the images (EoC)









It is again a question of features...

→ We speak about signatures





#### From CBIR to Word spotting



**Off-line indexation** 



#### From CBIR to Word spotting

#### **Results: retrieved images**

25. To Ensign Floming of the Virginia Regiment. You are hereby ordered to repair to baptain Hoggs Company at Fort Dinmedoic with eight good men as that bompany is without a Surgeon, if you will do that duty an 2700270.pnc allowance will be made you for it. you are to provide medicines, Ve. upon the best terms you can. This Order Sexpect will be immediately complied with; and that no Delays be offound ... you are to account with Captain Bell for your recruiting money, before you leave him. If you should arive at Augusta Court House 2710271.pnc before Sergeant Wilper and his Party you are to halt there until he joins, in order to escort the ammuneteon, Se. for the Fort, where you will receive blothes and arms for the men. 25. To Captain Bell of the Virginia Regiment. Thave Ordered Ensign Fleming

#### **Query images**



**On-line Retrieval** 



#### From CBIR to Word spotting





# Is it just a question of meta-data?



Dans la Mer d'Inde il y a une espèce de poissons qui ont

dans leur peau des poils si longs que les ger quant il les ont a en font des vêter pour s'habiller.











#### The standard model Descriptive meta-data + transcription

#### We have usually

- Descriptive meta-data in standard formats (MARC, EAD, Dublin Core, MODS, ...)
  - Edited manually
  - "Semantical" information
- Text transcription associated to additional meta-data (TEI)
  - Semi-automatic transcription or manually edited
  - "Semantical" information





#### The standard model Perceptual meta-data

- It seems that CBIR can help to extract and save supplementary information about image content (EoC) without going to the semantical aspect (recognition)
  - Regions of interest

  - □ Shapes, positions, colors, textures, … → Numerical values (vectors)





Is it just a problem of Meta-data? What about the encoding formats?

#### Structuration of data and file formats is a difficult problem

- Data architects are needed
- Some interesting formats linked with previous discussions
  - **METS Metadata Encoding and Transmission Standard**
  - $\circ~$  ALTO Analyzed Layout and Text Object
  - TEI Text Encoding Initiative





#### Is it just a problem of Meta-data? ALTO for OCR

#### ALTO = Analyzed Layout and Text Object

- Standard XML
- Created in 2003 during METAe project
- Developed by Graz, Linz, Innsbruck universities
- Description of the content and the physical layout of one page
- Used by several OCR software
- Adapted and used by the BNF and other libraries
- Drawbacks: huge / static





#### Is it just a problem of Meta-data? TEI for transcription and enriched contents

#### TEI = Text Encoding Initiative - Standard XML

- Content tagging and logical structure encoding (full document)
- Used a lot by libraries
- Too much « open »? → Quit « complex »





#### Is it just a problem of Meta-data? Link between meta data?

#### METS – <u>Metadata Encoding and Transmission Standard</u>

- Open XML Standard created in 2001 by the Digital Library Federation maintained by METS Editorial Board
- XSD-Schema
- Linking between multimedia objects
- Complete Description of digitized content (images, texts, audio, sculptures, ...)
- Physical / logical structures
- Descriptive Meta data (DC, MODS, MARC, ...)





# METS – Physical Structure





# METS – Logical Structure





# Conclusions



de poissons qui ont dans leur peau de

longs que les gen quant il les ont at en font des vêten



# Dans la Mer d'Inde il y a une espèce CONTENTS

ben a wheth a way





#### From Pixels to Contents Conclusions

#### Building tools for the valorisation of digitized historical content is a pluri-disciplinary task

- $\square$  Meta-data production  $\rightarrow$  Experts of the domains
- $\Box$  Selection and verification of the data  $\rightarrow$  Experts + Data accuratist
- $\square$  Structuration of the data and system  $\clubsuit$  Data / system architect
- Computer vision, Machine learning > Data scientist

#### Manual indexing is needed

- □ Descriptive meta-data → Semantical meta data
- Standard formats for data encoding
- □ Annotations could be seen as supplementary meta-data?

#### Operational methods and tools are available

- Acquisition devices
- □ Automatic tools: low level image processing, OCR
- Perceptual meta-data should be added : CBIR



#### From Pixels to contents Conclusions - Perspectives

#### The actual context: big data and heterogeneous collections

- □ Connexion between data, mutual enrichment, interoperability
- Introduction and management of additional knowledge
- Facing the diversity of the types of contents and usages

#### Quality of the interaction instead of only the quantity

Semantic Web : queries reformulation, smart crawlers, automatic categorisation





# A walkthrough around interactive systems



Dans la Mer d'Inde il y a une espèce de poissons qui ont dans leur peau de longs que les gen

quant il les ont at en font des vêten



## CONTENTS

1. The second second





#### A walkthrough around interactive systems AGORA





#### A walkthrough around interactive systems AGORA



□ EoC



#### AGORA

File       Project       Scenario       12 Operations       12 Op
Index: Parent-DocUMENT.New-CC.BnantzationBaseThreshold=ID           OCUMENT.New-CC.BnantzationBaseThreshold=ID           OCUMENT.New-CC.BnantzationBaseThreshold=ID
Background thread done. CCreateLabel: Create Label CC CCreateLabel: OK. Background thread started CC CC

## AGORA

#### 🖳 proj301.xml

File Project Scenario Misc. Operations L1 Operations L2 Operations L3 Operations View

CHAR{X=88,Y=1515,Width=10,Height=48}{w=1,[88.0] FindCC:Parent=DOCUMENT,New=CC,BinarizationBasicThreshold=150 CC{X=88,Y=1515,Width=10,Height=48}{w=1,[88.0 CreateLabel:Type=NOISE Classify:Type=ClassifyBySize,Parent=DOCUMENT,Child=CC,New=NOI CHAR{X=135,Y=1516,Width=25,Height=50}{w=1,[135 Delete:Parent=DOCUMENT,Child=NOISE,DeleteAll=True CC{X=135,Y=1516,Width=25,Height=50}{w=1,[13] CreateLabel:Type=CHAR Classify: Type=ClassifyBySize,Parent=DOCUMENT,Child=CC,New=CHA CHAR{X=69,Y=1513,Width=18,Height=53}{w=1,[69.0] CreateLabel:Type=IMAGE CC{X=69,Y=1513,Width=18,Height=53}{w=1,[69.0 Classify:Type=ClassifyBySize,Parent=DOCUMENT,Child=CC,New=IMA/ CreateLabel:Type=ACCENT CHAR{X=189,Y=1518,Width=25,Height=50}{w=1,[189] Classify:Type=ClassifyBySize,Parent=DOCUMENT,Child=CHAR,New=/ SetFeatures:Parent=DOCUMENT,Child=ACCENT,Mode={{a=Average/ CC{X=189,Y=1518,Width=25,Height=50}{w=1,[18: Insert:Parent=DOCUMENT,Childs={CHAR,ACCENT,},New=CHAR,Feat CHAR{X=367,Y=1518,Width=10,Height=50}{w=1,[367] CreateLabel:Type=LINE Fusion:Parent=DOCUMENT,Childs=CHAR,New=LINE,Comparator={0= CC{X=367,Y=1518,Width=10,Height=50}{w=1,[36] CHAR{X=380,Y=1519,Width=11,Height=48}{w=1,[380 CC{X=380,Y=1519,Width=11,Height=48}{w=1,[38] CHAR{X=907,Y=1524,Width=27,Height=50}{w=1,[907 CC{X=907,Y=1524,Width=27,Height=50}{w=1,[90] CHAR{X=818,Y=1528,Width=17,Height=29}{w=1,[818] CC{X=818,Y=1528,Width=17,Height=29}{w=1,[81] CHAR{X=762,Y=1529,Width=16,Height=49}{w=1,[762 CC{X=762,Y=1529,Width=16,Height=49}{w=1,[76; CHAR{X=794,Y=1529,Width=21,Height=28}{w=1,[794 CC{X=794,Y=1529,Width=21,Height=28}{w=1,[794 500 600 ..... Results = 635 Results = 635 Label Count CInsert: OK NONAME CInsert: LINE <- { LINE + LINE } CInsert: OK. DOCUMENT CFusion: Phase 3 : Simplification ... CC CDelete: Delete child LINE... NOISE 100 CHAR 200 300 IMAGE 400 ACCENT 500 LINE 600 700 800 900 1000 1100 1200 Results = 1270 CDelete: OK. CFusion: OK. Background thread done.

X=907 Y=705

PREMIER LIVRE faut laisser aus femmes céte vaine superstitio de parolles )auec vne espoge. Voila pourquoy, spongia, est vn motobscoene en Latin: & estoit céte esponge atachée au bout d'vn baston, comme tesmoignel'histoire de celuy qu'on menoit pour estre presenté aus bestes deuant le peuple, qui demanda congé d'aller a ses affaires, & la n'aiant autre moie de se tuer, il se fourra ce baston & espoge dans le gosier & s'en estoufa. Ilz s'essuyoint le catze de laine perfumée, quad ilz en auoint faich, As sili ail f. in falles and bestalan





#### AGORA

Proj301.xml	the second second second second second based based in a second second second second second second second second	
File Project Scenario Misc. Operations L1 Operations	L2 Operations L3 Operations View	
CreateLabel:Type=CC FindCC:Parent=D0CUMENT,New=CC,BinarizationBasicThreshold=150 CreateLabel:Type=NISE Classify:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=N01 Delete:Parent=D0CUMENT,Child=N0ISE,DeleteAll=True CreateLabel:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=CHA CreateLabel:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=IMAB Classify:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=IMA CreateLabel:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=IMA CreateLabel:Type=ClassifyBySize,Parent=D0CUMENT,Child=CC,New=IMA CreateLabel:Type=ClusSifyBySize,Parent=D0CUMENT,Child=CCENT,Mode={a-Average/ Inset:Parent=D0CUMENT,Child=CCENT,Mode={a-Average/ Inset:Parent=D0CUMENT,Child=CCENT,New=CHAR,Fea CreateLabel:Type=UINE Fusion:Parent=DINE,Childs=CHAR,New=LINE,Comparator={0= CreateLabel:Type=VORD Fusion:Parent=LINE,Childs=CHAR,New=WORD,Comparator={0= CreateLabel:Type=CHAR_New=WORD,Comparator={0= CreateLabel:Type=Childs=WORD,WORD,New=WORD,Coeff Delete:Parent=WORD,Child=WORD,DeleteAll=False	<ul> <li>WORD(X=760,Y=332,Width=128,Height=68}{w=5,[76</li> <li>CHAR(X=787,Y=332,Width=28,Height=49}{w=1,[7]</li> <li>CC(X=787,Y=332,Width=23,Height=20}{w=1,[7]</li> <li>CC(X=760,Y=352,Width=23,Height=20}{w=1,[7]</li> <li>CC(X=760,Y=352,Width=24,Height=29}{w=1,[7]</li> <li>CC(X=864,Y=352,Width=24,Height=29}{w=1,[7]</li> <li>CC(X=864,Y=352,Width=24,Height=29}{w=1,[7]</li> <li>CC(X=864,Y=353,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=353,Width=28,Height=29}{w=1,[7]</li> <li>CC(X=801,Y=866,Width=18,Height=30}{w=1,[858]</li> <li>CC(X=858,Y=866,Width=18,Height=30}{w=1,[858]</li> <li>CC(X=858,Y=866,Width=18,Height=30}{w=1,[6]</li> <li>CC(X=858,Y=866,Width=18,Height=30}{w=1,[6]</li> <li>WORD(X=451,Y=888,Width=45,Height=45}{w=1,451}</li></ul>	ot ob-
CFusion: Phase 3 : Simplification CDelete: Delete child WORD 100 200 300 400 500 600	WORD(V=43); 1=203; Vitath=43); W=1; [4:5]       CHAR(X=451; Y=858; Width=45, Height=45); W=1; [4:5]       CC(X=511; Y=858; Width=45, Height=45); W=1; [5:8]       WORD(X=508; Y=871; Width=17, Height=30); W=1; [5:8]       CHAR(X=508; Y=871; Width=17, Height=30); W=1; [5:8]       CC(X=5108; Y=871; Width=17, Height=30); W=1; [5:8]       COULT       NONAME       DOCUMENT       CC	omme
700 800 Results = 878 CDelete: OK. CFusion: OK. Background thread done. Background thread started CIntersect: ChildsEOC -> NewEOC CInsert: WORD <- { WORD + WORD } Results = 33 Results = 33	NOISE CHAR IMAGE ACCENT LINE WORD	on me-
CInsert: OK. CIntersect: OK. Background thread done. CDelete: Delete child WORD Results = 66 CDelete: OK.	e aus bei	tes de-



#### AGORA

#### proj301.xml

56

X=764 Y=335

- 0 ×





#### AGORA





- 0 -X

44 37

275,714

443,400

#### A walkthrough around interactive systems **Retro: Typographic analysis**

Body Height Measure

ryphius.

rvlict, 2007, no. 47.

D:\Dataset\VERVLIET\_Conspectus\CESR\_56656\_0089.jpg

CS.

B

Select Frontier

[20] (in mm)

[x] (in mm)

Consider Line Space

Top (in pixel)

Top (in pixel)

Bottom (in pixel)

200

3

, Paris, G. Bossozel (Evangelium; Morea

HIERONYMVS IN CA.

talogo fcriptorum eccle=

fiasticorum. MATTHAEVS qui & Leui,

ex publicano apoitolus, primus in lu

propter eos qui ex circuncifione credide

rant, euangeliu Chrifti Hebraicis literis verbif q compoluit.Quod quis postea in

Græcum transtulerit, no fatis certum eft.

Porro ipfum Hebraicum habetur víque hodie in Cafariensi bibliotheca, quam,

Paphilus martyr fludiofiffime confecit.

Mihi quo'g à Nazaræis , qui in Bœrea

(Mignonne); 20 44 XI : 1.5.

#### Specific GUI

- Font models creation
- Model visualization
- Fonts Measurement and identification
  - 1. Select a page
  - 2. Selection of the frontier (top/bottom for [20], [x], [1]height), and display of measured values
  - Estimated Body height designation
  - Possibility to export computed information in xml





#### A walkthrough around interactive systems Retro: Typographic analysis

Data generation

59

- Which formats (XML + binary, color images)
- □ Which meta-data for characters, glyphs, patterns (unicode, ...)
- Training data generation
  - Which fonts ? Which models ? How many ?
  - Clean ? Degraded ?
- Early modern font classification
- Estheatic study of fonts
- Font Base Batyr du CESR
- Available data on Paradiit web site
  - The 'Gering' Pica-Roman [R80]
  - Cicéro (1478) (cf. [Vervliet2010] N°55)
  - Garamont's Great Primer Roman [R118]
  - □ Gros-romain (1549) (cf. [Vervliet2010] N°119)
  - Vérard Gothic Bâtarde Great Primer Roman



R P K & P R

A walkthrough around interactive systems Retro: Assisted Transcription

#### EoC Clustering → Assisted transcription / indexing

- Feature selection
  - Pixels (NdG ? B&W ?)
  - Visual Features

#### Comparison algorithms and metrics

- Millions of matching
- Cluster representatives
- Time complexity
  - □ 1 books = several days → distribution
  - □ M. de files → Mass of data



RETRO 2012

ile Page Clustering Transcription Results Typography Hel

fermer il fe fourra ce hafton & efo

tibi nil faciam, fed lota metula la v auoit aus carrefoursa Rome

aiffeaus & demy cu piffer aus paffans.



# These tools are still to improve...



## **Extraction & Clustering of EoC**

#### Montaigne - 1557

- 119 pages, 3260 blocs de texte centraux
- 125 744 composantes connexes (pseudo caracteres)
- 29 943 clusters
- 25 classes = 25% du texte
- 136 classes = 50%
- 1500 classes = 70%
- 20 000 classes = 90%
- 79% des classes comportent un seul élément
- Classe la plus grosse = 3%

1% des formes sont mal classées

Omnem crede diem tibi diluxisse jupremum Grata superueniet quæ non sperabitur bora.

Il eft incertain ou la mort nous atten de, attendons la par tout. La premeditation de la mort eft premeditation de la liberté. Qui a apris a mourir ila defapris a feruir. Le fçauoir mourir nous afranchit de toute fubiection & contrainte. Paulus Æmilius refpondit a celuy que ce miferable roy de Macedoine



## **Extraction & Clustering of EoC**

#### Vésale – Latin - 1543

- 150 pages avec notes en marges
- 1.062.081 composantes connexes (pseudo caractères)
- 40.000 classes (clusters)
- 57% des classes comportent un seul élément
- 90% des classes comportent moins de 10 occurrences
- Les ignorer durant la transcription correspond à un caractères manquants sur 14 → intolérable !!!





Pourquoi ?

- Bruits, dégradations
- Caractères connectés
- Caractères cassés

0 ...

#### DE HVMANI CORPORIS FABRICA LIBER 1: 43 diftinguente, alia educitur futura, oblique deorfum inter caninum & inciforium canino proxi

mum delata, quæ cum illa communis efficitur, quæ in palati extremo iuxta inciforios dentes tranfuerfim in illis animalibus fertur. Hæc futura citra omnem cartilaginis interuentum adeò infignis eft, ut quartum à me enumeratum maxillæ os, canes in duo diulium commonfirent. Verum non arbitrandum eft futuram hanc, inter dentium præfipiola in canibus aut fimijs duetam, à fuperciliorum medio deorfum continuo progreffu (& fi Galenus ita doceat) ferri, fed uti làm admonuimus, ex medio ferè ducfu illius futuræ incipit, quæ externum nafi offium latus terminat. Quod ubi in cane fimia ue infpexeris, accurate expendito, quàm impoffibile fit, hanc futuram porrigi inter dentes hominis, maxillam adepti breuiffimam, & dentibus caninis prorfus exiguis donati. Infuper animaduertito, quid in libris de Partium ufu Galenus fibi uoluerit, quum duodecim fuperioris maxillæ offa effe obiter feripfittut & Introductorij feu Medici autor, duodecim tantum maxillæ offa effe obiter feripfittut & Introductorij feu Melenus hæc offa prolisius comemorans, nouem tanti recenfer, nullumé os priuatim inciforijs emo.

> DE MAXILLA INFERIORI, Caput X.



VTRARVMQVE DECIMI CAPITIS FIGVRARVM, ET earundem characterum Index. PRIOR feu dextra buius capitis figura, inferiorem maxillam und cum dentibus an



#### Try Agora and Retro (Practical works) $\rightarrow$

http://rfai.li.univ-tours.fr/PagesPerso/jyramel/fr/cours\_PEEN.html

#### Project web site →

https://sites.google.com/site/paradiitproject/

