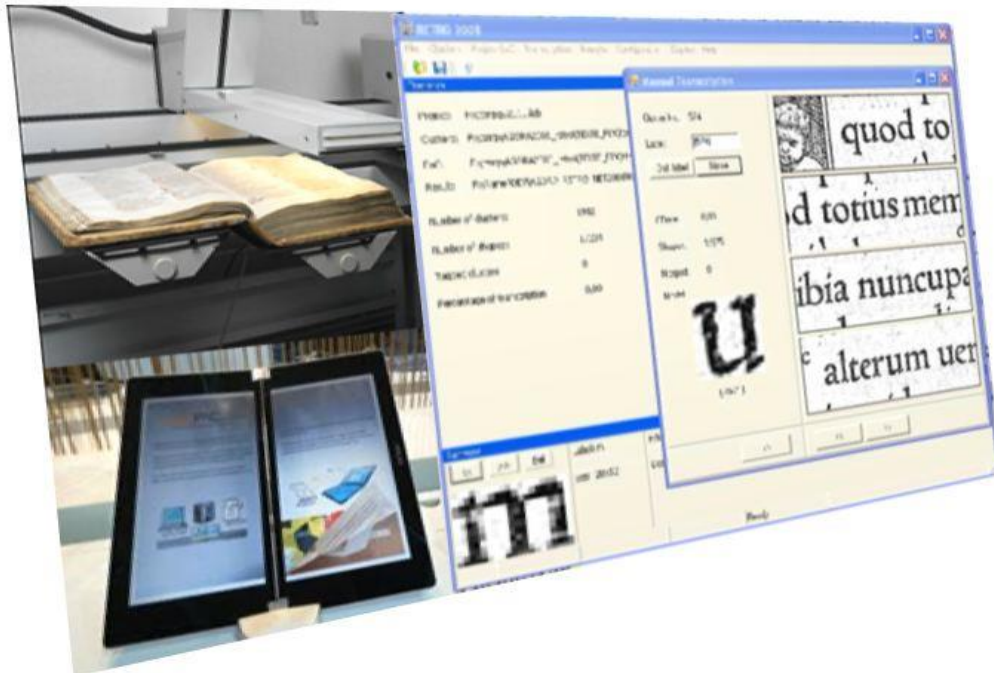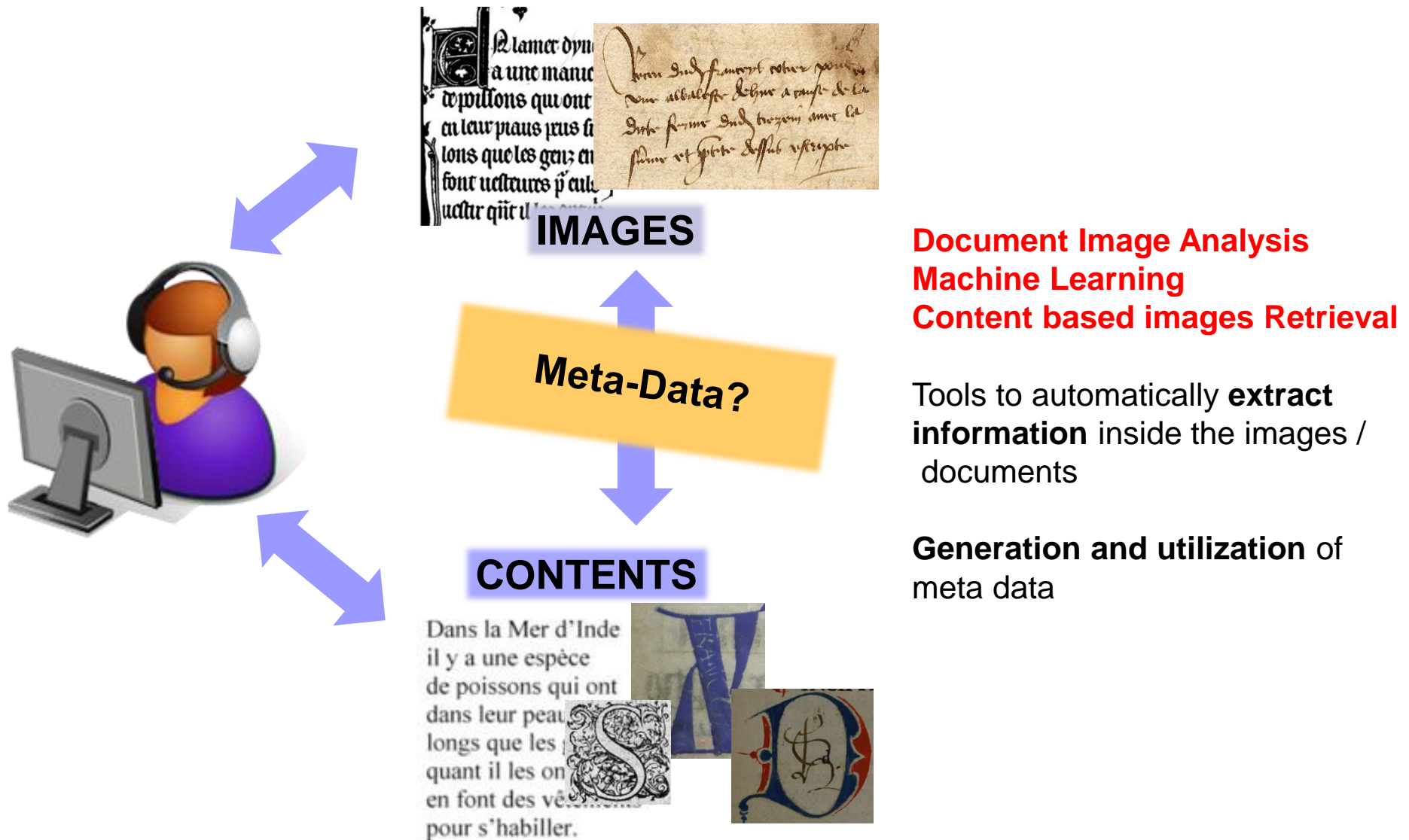# From pixels to content: An overview of the main techniques used in DIA

**Jean-Yves RAMEL**

# From Pixels to Contents
## Introduction
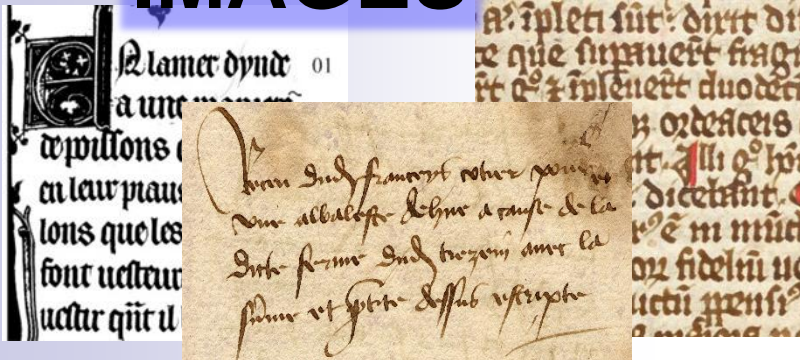
**IMAGES**

**Meta-Data?**

**CONTENTS**

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

**Document Image Analysis**
**Machine Learning**
**Content based images Retrieval**

Tools to automatically **extract information** inside the images / documents

**Generation and utilization** of meta data

# 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?
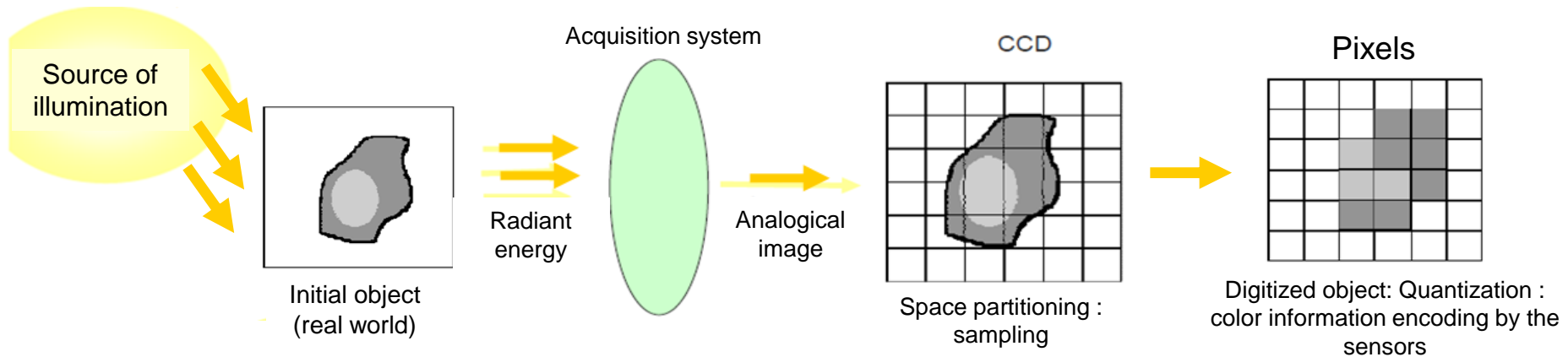
- **Conclusions and perspectives**

# From Pixels…

**IMAGES**

**From Pixels…**
# Digitization ➜ Set of pixels ?

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



Source of illumination

Acquisition system

CCD

Pixels

Radiant energy

Analogical image

Initial object (real world)

Space partitioning : sampling

Digitized object: Quantization : color information encoding by the sensors

- **Quantization**
  - ☐ Assignment of a numerical value drawn from the received lighting energy / pixel (grid unit)
  - ☐ Continuous value (xi,yi) → Discrete value (xi,yi) → Pixels
  - ☐ The range of colors that each pixel can take

# From Pixels…
# What is an image?

## Image Quantization

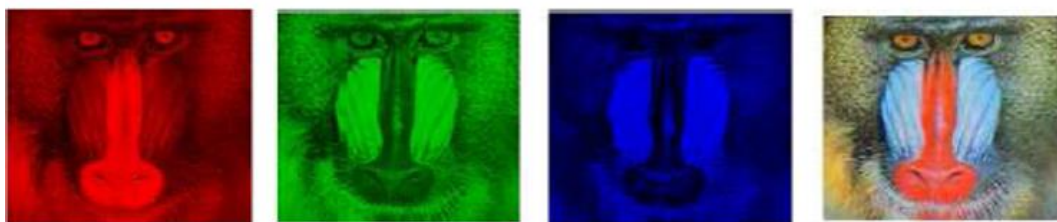**Binary images:** *I(i,j)* = 0 black    or    *I(i,j)* = 1 white

**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

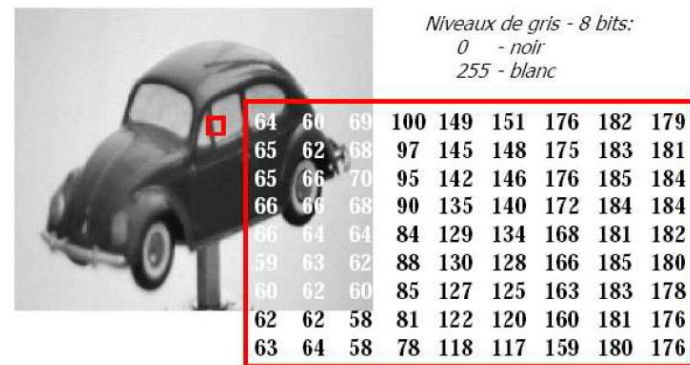$I_1(i,j)$ = 0…..255   –   $I_2(i,j)$ = 0…..255   –   $I_3(i,j)$ = 0…..255

## 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 *I(i,j)* associated to each pixel *s(i,j)* represent their brightness intensity


32 bits          6 bits          1 bit



Niveaux de gris - 8 bits:
0   - noir
255 - blanc

| 64 | 60 | 69 | 100 | 149 | 151 | 176 | 182 | 179 |
|----|----|----|-----|-----|-----|-----|-----|-----|
| 65 | 62 | 68 | 97  | 145 | 148 | 175 | 183 | 181 |
| 65 | 66 | 70 | 95  | 142 | 146 | 176 | 185 | 184 |
| 66 | 66 | 68 | 90  | 135 | 140 | 172 | 184 | 184 |
| 66 | 64 | 64 | 84  | 129 | 134 | 168 | 181 | 182 |
| 59 | 63 | 62 | 88  | 130 | 128 | 166 | 185 | 180 |
| 60 | 62 | 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



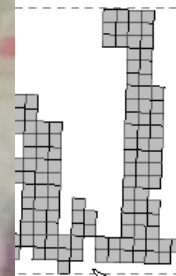256x256    128x128    64x64    32x32

- **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?

...olated patterns can correspond ...pixels) !

...on the boundaries of the shapes ...ocess

...arance of characters until the

...acters or touching characters
**...e resolution**

Touching

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, …

Corriger l'image (rotation, wrapping)

**From Pixels…**
# Image (pre-)processing

Curvature and skew correction is possible on text images

point & le reste de satin blanc, & tout passementé & pou
filé d'or : celluy des Espingliers bonnet, collet, chausses,
souliers de uelours noir:le pourpoint de satin cramoisy,
doubleure des chausses correspondât, rayez de passemen
& traisses d'or. Apres lesquelz passoient quelques premie
rangs armez & accompaignez de deux centz & sept Tif
rans portantz rouge & noir:les troys Enseignes derrie
eulx braues & bien en ordre, & marchantz deuât deux cer
cinquâte six Cordoanniers uestus de blanc & noir, laissar
à leurs espaules les troys Lieutenantz autant brauement
ordre, & conduisantz centz quatre uingtz & douze Esp:
gliers portantz le pourpoint de uelours, satin, ou taffe
rouge, le collet & bonnet noir auec plume blanche, & gr:
satissaisant à chascun.

Tout d'un ordre suruint la sixiesme Bande autant be!
que plaisante pour la diuersité des couleurs:laquelle côm
ca par le rang de ses troys Capitaines de Rue neuue acc
itré de uelours noir, blanc, & bleu mouchetté menuem
de bouttons d'or, accôpaigné du Capitaine dès Chappeli
uestu de uelours blanc & noir & uerd à petitz grains d'
suyuant d'un mesme pas auec celluy des Fondeurs en ha
de uelours blanc, & noir, & aurangé, recamé & bisetté d
gent. Et lequel rang auec ses Tabourins & Fiffres de mes
fut suyuy d'aucuns autres armez de corseletz & animes; &
suytre de Rueneuue en liuree de noir blanc & bleu, &
nombre de quatre centz uingt & troys:lesquelz estoient
stez de troys Enseignes suyuantz auec mesmes couleur:
leurs enseignes, guidantz apres eulx cent soisante & f
Chappellier de blanc noir & uert:Et à la file les troys Li

point & le reste de satin blanc, & tout passementé & pou
file d'or : celluy des Espingliers bonnet, collet, chausses,
souliers de uelours noir:le pourpoint de satin cramoisy,
doubleure des chausses correspondât, rayez de passemen
& traisses d'or. Apres lesquelz passoient quelques premie
rangs armez & accompaignez de deux centz & sept Tif
rans portantz rouge & noir:les troys Enseignes derrie
eulx braues & bien en ordre, & marchantz deuât deuv cer
cinquâte six Cordoanniers uestus de blanc & noir, laissar
à leurs espaules les troys Lieutenantz autant brauement
ordre, & conduisantz centz quatre uingtz & douze Esp:
gliers portantz le pourpoint de uelours, satin, ou taffe
rouge, le collet & bonnet noir auec plume blanche, & gr:
satissaisant à chascun.

Tout d'un ordre suruint la sixiesme Bande autant be!
que plaisante pour la diuersité des couleurs:laquelle côm
ca par le rang de ses troys Capitaines de Rue neuue acc
itré de uelours noir, blanc, & bleu mouchetté menuem
de bouttons d'or, accôpaigné du Capitaine des Chappeli
uestu de uelours blanc & noir & uerd à petitz grains d'
suyuant d'un mesne pas auec celluy des Fondeurs en ha
de uelours blanc, & noir, & aurangé, recamé & bisetté d
gent. Et lequel rang auec ses Tabourins & Fiffres de mes
fut suyuy d'aucuns autres armez de corseletz & animes, &
suytre de Rueneuue en liuree de noir blanc & bleu, &
nombre de quatre centz uingt & troys:lesquelz estoient
stez de troys Enseignes suyuantz auec mesines couleur:
leurs enseignes, guidantz apres eulx cent soisante & f
Chappellier de blanc noir & uert:Et à la file les troys Li

**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

**IMAGES** → **CONTENTS**

**From pixels … to text**
# Automatic transcription but also layout Analysis



1-Offrant appres plusieurs rues sur ce faites
2-au lieu et temps acostume aly livre et
3- accenste la somme et quantite dessoulz
4-estempte du voloir et consentement de
5-plusieurs des borgeys de la dite ville
6-de chastellion                                    _____IIIIxXX  V
flor(ins)

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-                                    _____II flor(ins)

12-receu de Alexandre ? escoffier al loup ?
13- dudit lieu de chatellion pour la
14-rense et ferme du commun de la dite
15-ville et communaute ? pour ung an commettre ?
16-alaleste nativite fi jehan baptiste lan
17-mil quatre cent et cinquante? A ly par le
18-pris dessouls esempt accense et hure ?
19-ce 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 ?

**Preservation of the link between the Text and the Image**

**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

- **Color Image ➔ Grayscale ➔ binarisation**



**EoC**

**From pixels … to text**
# An overview of OCR mechanisms

## Just to illustrate the difficulties…

■ Most of the segmentation methods need a binarisation



**Threshold selection ?**

**From pixels … to text**
# An overview of OCR mechanisms

## Just to illustrate the difficulties…

- 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

**From pixels … to text**
# An overview of OCR mechanisms

## First step: Image segmentation / Connected components

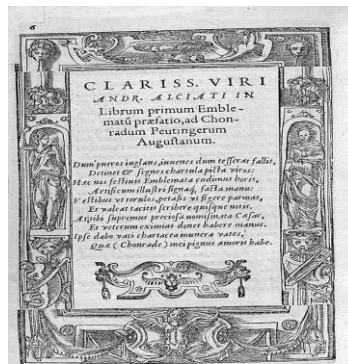- Then, we can try to group black pixels together to **localize** and **recognize** higher level Element of Content **(EoC)**

**From pixels … to text**
# An overview of OCR mechanisms

## Next step : Layout analysis

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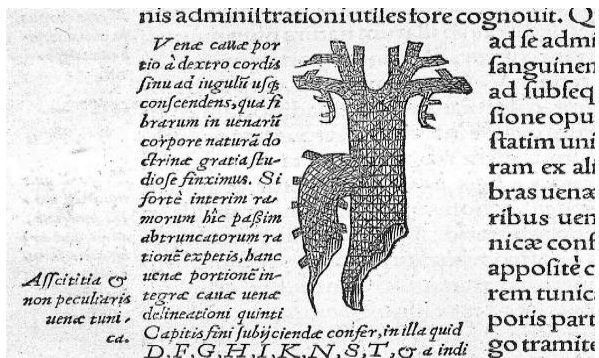- Choosing how to organize the XML tree (physical / logical) is not so easy…

**From pixels … to text**
# An overview of OCR mechanisms

## 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)



**Physical structure ➔**        **← logical structure**

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

## 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



Projet BVH - Paradiit (https://sites.google.com/site/paradiitproject/)➔
AGORA: an user-driven system for content extraction in historical printed books

Projet Européen Meta-e (http://meta-e.aib.uni-linz.ac.at/ )➔
First commercial system for automatic layout analysis Dutch books of XVIIIe

**From pixels … to text**
# An overview of OCR mechanisms

**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)



V=(0,3,0,4,12,4,3,0,3)    V=(6,10,0,12,4,0,10,10,0)

A          E

… [ Training set ] …

1 EoC ➔ 1 Vecteur $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ .. \\ x_n \end{bmatrix} \in I\!R^n$

[ 2D Representation of the training set]

Stable and discriminative features

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

## 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)



$V=(0,10,0,5,14,5,3,0,2)$

$V=(0,3,0,4,12,4,3,0,3)$          $V=(6,10,0,12,4,0,10,10,0)$

?          Forme inconnue          ?

$V=(0,10,0,5,14,5,3,0,2)$

[ 2D Representation of the training set]

$$D(A,?) = \sqrt{(0-0)^2 + (3-10)^2 + (0-0)^2 + (4-5)^2 + \ldots + (3-2)^2}$$

$$D(A,?) = 7,48 \text{ et } D(B,?) = 19,05$$

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

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

**From pixels … to text**
# An overview of OCR mechanisms

- **Why commercial OCR are not working well on historical documents?**

  - Noises and degradations
  - Unusual layout
  - Unsuited training set

**Fine Reader**

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

- **Why commercial OCR are not working well on historical documents?**

- Lack of data, knowledge and experiences
  - ☐ Unusual fonts and characters ➔ 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

**IMAGES**

**CONTENTS**

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



- **Ornamental letters ( +of 20000)**



- **Figures (+ de 1500)**



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

**From pixels … to non-textual contents**
# An overview of Content Based Image Retrieval

# Perceptual meta data instead of classical meta data

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

- Computation time not crucial
- **Signatures ➜ Visual Features**
  **➜ Perceptual meta data**

**Set of images to index**

$(x_0^0, x_1^0, \ldots, x_p^0)$

$(x_0^n, x_1^n, \ldots, x_p^n)$

*Base de données images*

Index
database

$\{M, C, \ldots\}$

Statisrics

# From pixels … to non-textual contents
# An overview of Content Based Image Retrieval



OFF-LINE **Back-office**

ON-LINE **Front-office**

**Using images as query instead of words**

DATABASE

Indexing

Image  3D-object  Video

**Search Engine**

Search

Relevance feedback

Classification

User interaction

QUERY

LABELS

RESULTS

DESCRIPTORS INDEX

**From pixels … to non-textual contents**
# An overview of Content Based Image Retrieval

It is again a question of features…

➔ We speak about signatures

**Query Image**

**Signature computation**

**On-line retrieval**

- - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Off-line indexation**

**Index database containing perceptual signatures**

**Image Dataset**

**similar Images**

*Zone de similarité*

**Representation space**

**From pixels … to non-textual contents**
# An overview of Content Based Image Retrieval

- ## From CBIR to Word spotting



**Off-line indexation**

Document

EoC

Extraction

Signature computation

Clustering

Index

Indexing

>that (p.12, p.45, …)

>the (p.12, p.34, …)

the

that

From pixels … to non-textual contents
# An overview of Content Based Image Retrieval

■ From CBIR to Word spotting

**Query images**

**Results: retrieved images**

**On-line Retrieval**

**From pixels … to non-textual contents**
# An overview of Content Based Image Retrieval

## ■ From CBIR to Word spotting

# Is it just a question of meta-data?

**IMAGES**

**CONTENTS**

metadata

Dans la Mer d'Inde
il y a une espèce
de poissons qui ont
dans leur peau des poils si
longs que les ge[...]
quant il les ont a[...]
en font des vête[...]
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



**"Semantical" meta-data**

# 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
  - Visual features ➔ Perceptual **signatures and index**
  - Shapes, positions, colors, textures, … ➔ Numerical values (vectors)



**+ "Perceptual" meta-data**

**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
  - o **METS – Metadata Encoding and Transmission Standard**
  - o ALTO – Analyzed Layout and Text Object
  - o 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 »

```
<TEI>
    <text>
        <body>
            <table rows="5" cols="3"/>
            <p xml:id="par0"/>
            <p xml:id="par1"/>
            <p xml:id="par2"/>
            <p xml:id="par3"/>
            <div n="3" type="partie"/>
            <figure/>
            <div n="4" type="partie"/>
        </body>
    </text>
</TEI>
```

| Tableau | Figure |
|---|---|
| Paragraphe | |
| Paragraphe | Partie 4.0 |
| Paragraphe | |
| Paragraphe | |
| Partie 3.0 | |

Numéro de page

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

## METS – <u>M</u>etadata <u>E</u>ncoding and <u>T</u>ransmission <u>S</u>tandard
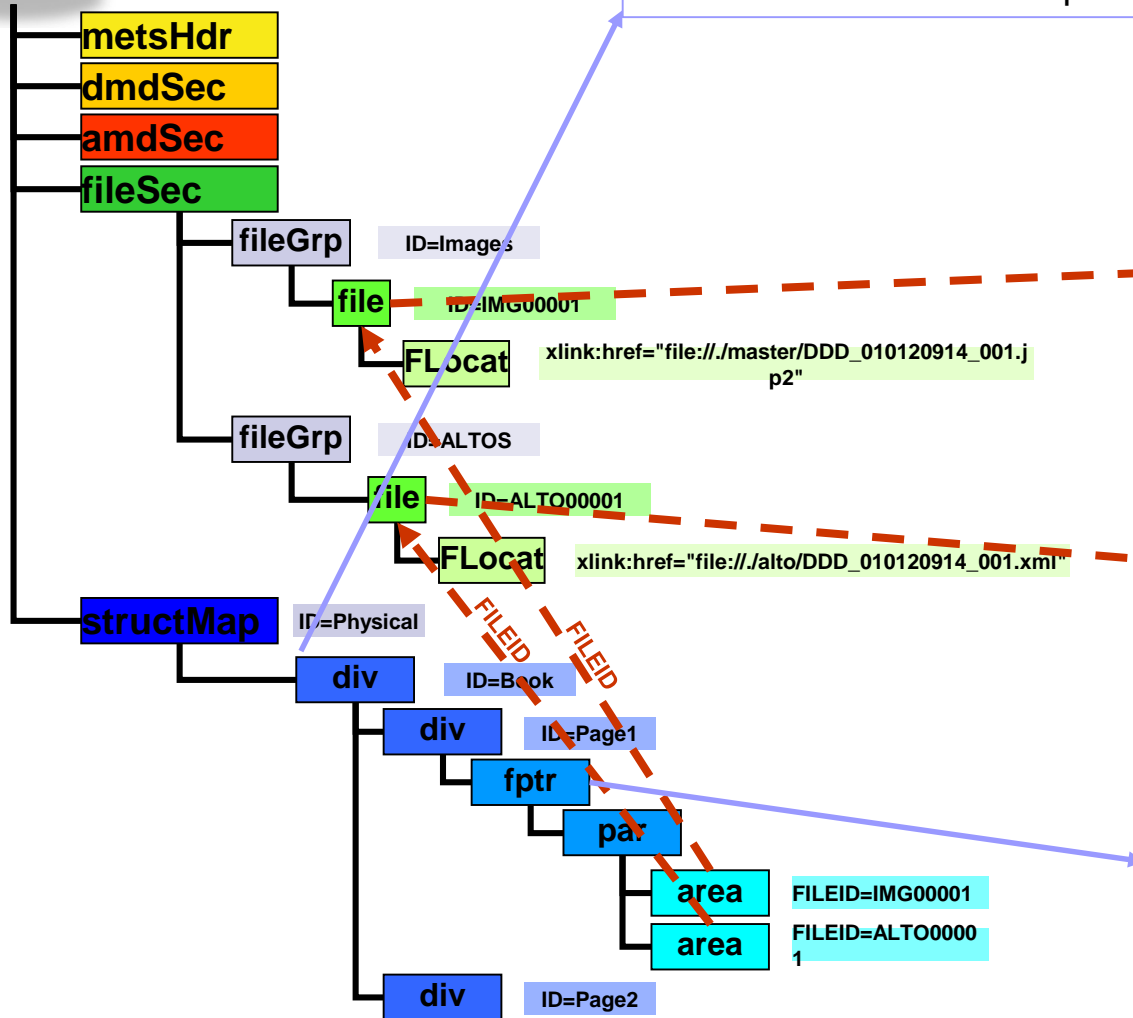
- Open XML Standard created in 2001 by the Digital Library Federation maintained by METS Editorial Board

- XSD-Schema

- **Linking between multimedia objects**

- Complete Description of digitized content (images, texts, audio, sculptures, …)

- Physical / logical structures
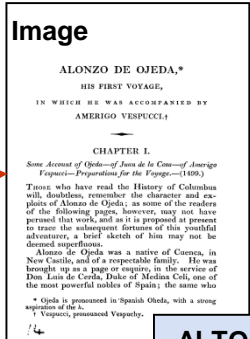
- Descriptive Meta data (DC, MODS, MARC, …)

- …

## METS – First level Elements

**metsHdr** — METS document header (info author, sotware, …)

**dmdSec** — descriptive metadata section (bibliographic notice)

**amdSec** — administrative metadata section (copyright, …)

**fileSec** — file inventory section (file localisation)

**structLink** — structural map linking (link between structures)

**structMap** — structural maps (physical et logical)

# METS – Physical Structure

METS

- metsHdr
- dmdSec
- amdSec
- fileSec
  - fileGrp — ID=Images
    - file — ID=IMG00001
    - FLocat — xlink:href="file://./master/DDD_010120914_001.jp2"
  - fileGrp — ID=ALTOS
    - file — ID=ALTO00001
    - FLocat — xlink:href="file://./alto/DDD_010120914_001.xml"
- structMap — ID=Physical
  - div — ID=Book
    - div — ID=Page1
      - fptr
        - par
          - area — FILEID=IMG00001
          - area — FILEID=ALTO00001
    - div — ID=Page2

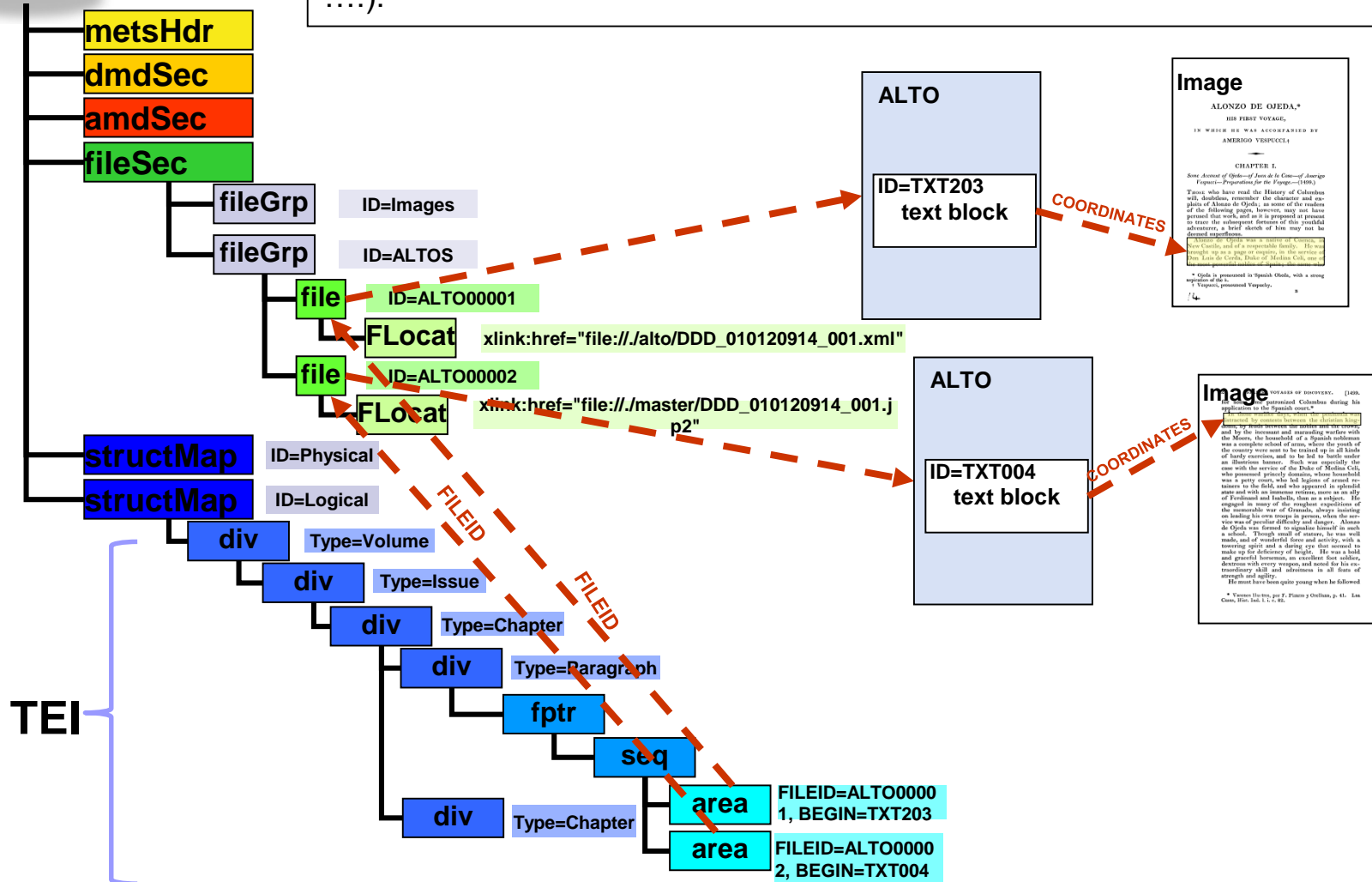Possible link to the descriptive meta data (dmd Sec)

METS allows to specify the locations of resource files ➔ File pointers METS (fptr)
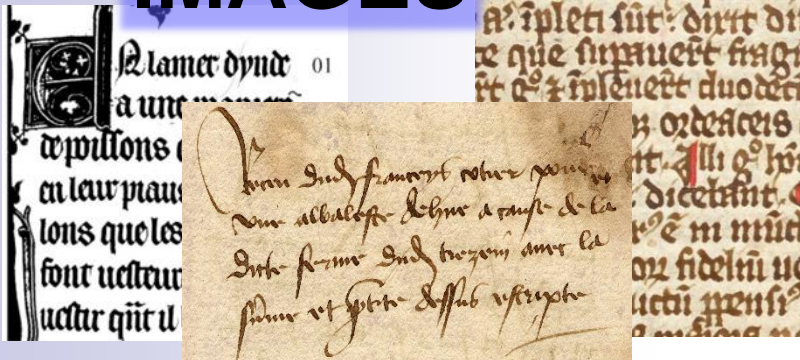
FILEID

# METS – Logical Structure



The Logical StructMap reflects the enrichment of the different logical blocks located in different pages that can be split by other logical blocks (like foot-notes, ….).
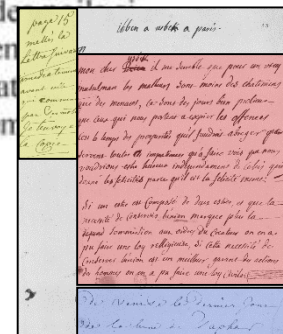
**METS**
- metsHdr
- dmdSec
- amdSec
- fileSec
  - fileGrp — ID=Images
  - fileGrp — ID=ALTOS
    - file — ID=ALTO00001
      - FLocat — xlink:href="file://./alto/DDD_010120914_001.xml"
    - file — ID=ALTO00002
      - FLocat — xlink:href="file://./master/DDD_010120914_001.jp2"
- structMap — ID=Physical
- structMap — ID=Logical
  - div — Type=Volume
    - div — Type=Issue
      - div — Type=Chapter
        - div — Type=Paragraph
          - fptr
            - seq
              - area — FILEID=ALTO00001, BEGIN=TXT203
              - area — FILEID=ALTO00002, BEGIN=TXT004
        - div — Type=Chapter

**TEI**

**ALTO**
ID=TXT203
text block

COORDINATES

**Image**
ALONZO DE OJEDA,*
HIS FIRST VOYAGE,
IN WHICH HE WAS ACCOMPANIED BY
AMERIGO VESPUCCI.

**ALTO**
ID=TXT004
text block

COORDINATES

**Image**

FILEID

FILEID

# Conclusions

**IMAGES** → **CONTENTS**

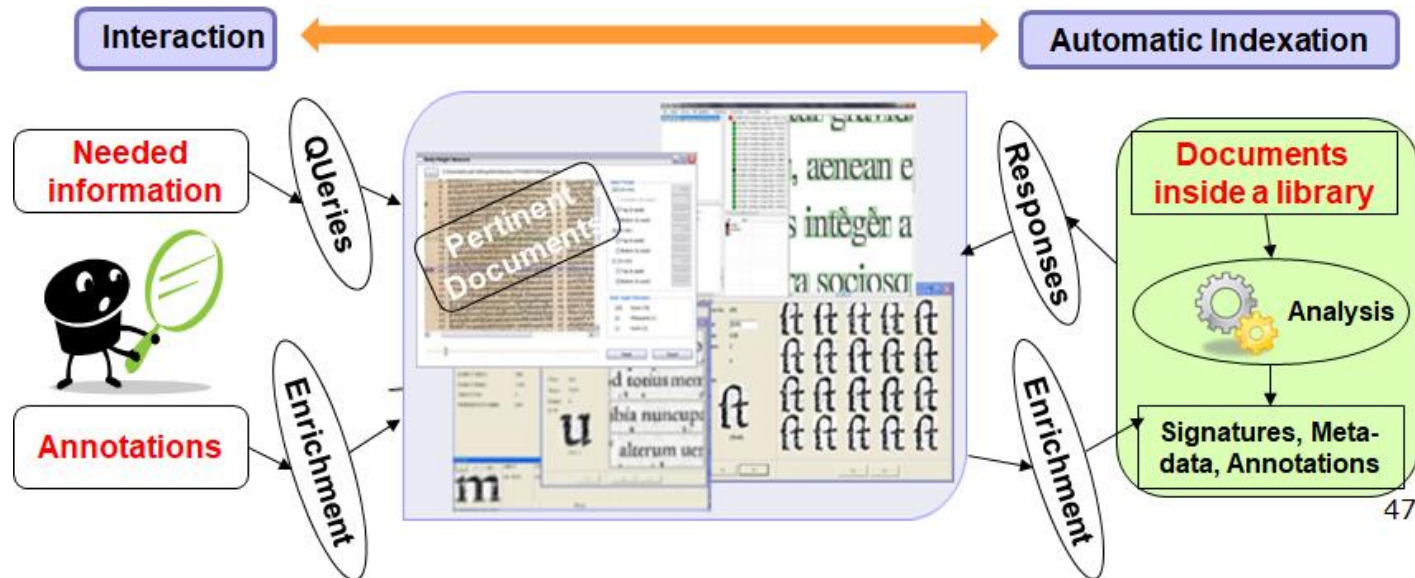**From Pixels to Contents**
# Conclusions

- **Building tools for the valorisation of digitized historical content is a pluri-disciplinary task**
  - ☐ Meta-data production ➜ Experts of the domains
  - ☐ Selection and verification of the data ➜ Experts + Data accuratist
  - ☐ Structuration of the data and system ➜ 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



47

# Thanks...

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