

Les Images de Documents

*Numérisation, Dématérialisation, Analyse,
Exploitation*

Véronique EGLIN

LIRIS, INSA Lyon

Janvier, 2011



Objectifs du cours

.Comprendre ce qu'est le document numérique:

- Pourquoi numériser ?
- A quoi servent les images de documents?
- Comment les numériser ?
- Et après que faire avec ?
- ... Ensuite les choses se compliquent: il faut les organiser, les indexer, les exploiter, les analyser...

.Quelques applications de la gestion documentaire dans l'industrie (documents d'entreprises), dans les bibliothèques (documents du patrimoine).

Plan du cours

Jeudi 6 janvier

- Présentation générale: les enjeux de la dématérialisation et de la GED
- Les bases du prétraitement des images (partie 1)

Vendredi 7 janvier

- Les bases du (pré-)traitement des images et de l'indexation par le contenu des documents multimédias (partie 2)

Jeudi 13 janvier

- Les bases de la Reconnaissance de formes et de la classification

Vendredi 14 janvier

- L'analyse de structures: physique et logique
- Des applications cibles: WS, CBIR, TAO, Tri, e-Edition...

Contexte

•A l'origine: des besoins d'accéder à des ensembles de documents

- Archives
- Documents juridiques, administratifs, ...
- Documents du patrimoine

•Généralement, on dispose de collections hétérogènes

- Différentes langues
- Différentes structures/mises en page
- Différentes provenances
- Différentes modalités d'acquisition et qualités...

•Rarement, on dispose de métadonnées descriptives des contenues, homogènes et pertinentes

Du document “papier” au document électronique?

.Objectif: “sans papier”

- Transfert efficace
- Organisation facilitée
- Reproduction “à la demande”...

.Accès à une plus grande variété de contenus

- Echange de données (email, attachments, diffusion...)
- Lorsque la “copie” rend un service identique à l’original

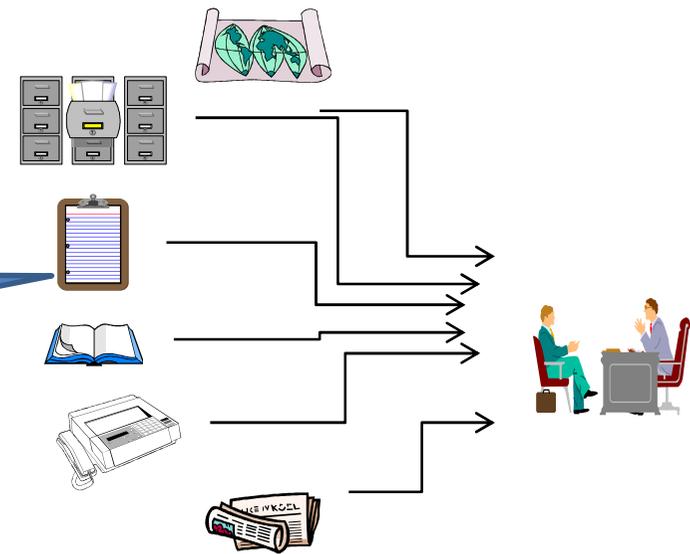
.Sauvegarde et préservation

.Edition électronique, édition critique

Quels usages?

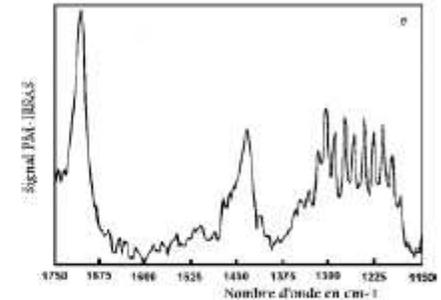
- .Internet
- .Email Attachments
- .Online Proceedings
- .Electronic Fax
- .Mass Digitization Repositories...
- .Sauvegarde du patrimoine

Une nécessaire organisation GED des documents



Quels usages?

- Médical : détection de pathologie à partir d'ECG/EEG, imagerie : aide à la décision
- Militaire : détection de cibles
- Sécurité automobile : détection d'obstacles
- Traitement du signal : reconnaissance de la parole, de l'écriture
- etc.

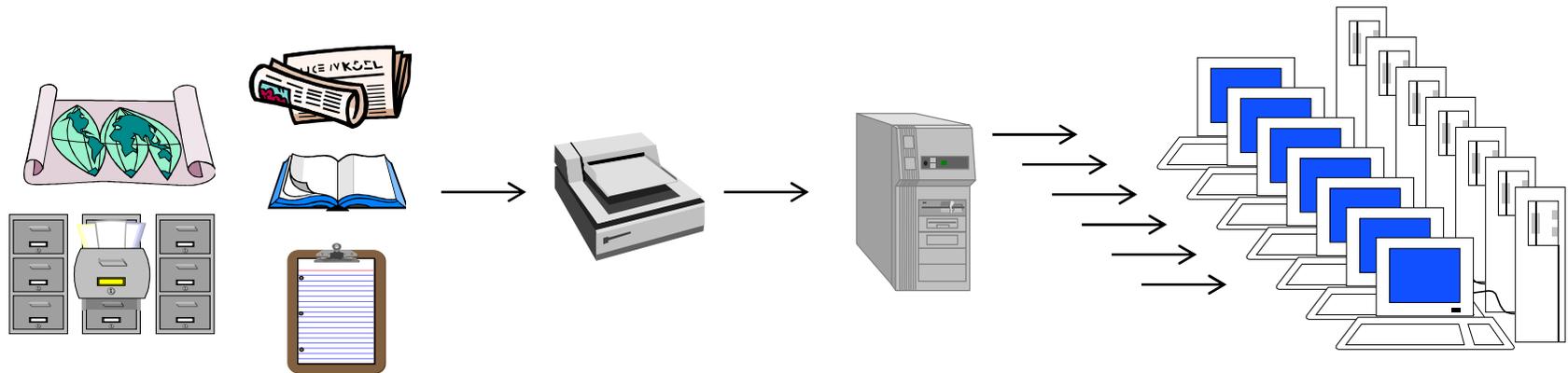


*Odria 7 ou 8 lithographies de
Balzac, elles sont petites
J'ai fait photographies un très
beau pastel (Matisse) de Gust
qui est au musée de Cour, où
il y a aussi un dessin de Boulangot.*



Bases d'images de documents

- .Collection d'images numérisées
- .Organisées pour rendre possible différentes actions tq: indexation, recherche, édition, dissémination, interprétation....



Définitions générales

Numérisation : *Conversion des documents sous la forme d'une image numérique à l'aide d'un scanner ; c'est la première étape de la chaîne de dématérialisation.*

OCR : *Reconnaissance Optique des Caractères, étape intermédiaire vers la dématérialisation; elle produit un texte ASCII au kilomètre, sans aucune structure, enrichi seulement des styles et des polices de caractères.*

ARD : *La Reconnaissance Automatique des Documents produit un document électronique enrichi de la structure logique et physique*

Dématérialisation ou Rétroconversion : *La dématérialisation consiste à convertir les documents papier dans un format logique pivot les rendant de nouveau rééditables et donc réutilisables. Elle va au-delà de la simple conversion ASCII du texte par un OCR, et propose une réelle réhabilitation de la structure interne du document, la plus proche possible de celle qu'il avait à l'origine. On parle dans ce cas de rétroconversion [Belaid 97]*

Chaîne de traitement de l'ARD

1/ Numérisation

2/ Pré-traitement

prétraitements →

- Binarisation
- Débruitage
- Correction d'inclinaison
- etc.

3/ Extraction d'attributs et indexation par le contenu

Extraction de caractéristiques →

- Augmenter la pertinence de la représentation
- Réduction de la quantité d'info.

4/ Classification, apprentissage et reconnaissance

Apprentissage automatique →

Idée = Apprendre la diversité et la variabilité des formes grâce en soumettant à la machine de (très) nombreux exemples

“machine learning”



Les tâches usuelles en ARD

- Analyse et reconnaissance des structures physique/logique
- Classification de zones d'intérêt (ROI)
- Recherche d'information (WS, Shape Spotting, CBIR...)
- Indexation et Reconnaissance des contenus des ROIs (OCR, transcription, identification de scripteurs, authentification....)
- Tri et applications industrielles

Partie 1

LA NUMERISATION



La révolution du numérique

La numérisation ne fait qu'une copie NUMERIQUE de l'ouvrage.

C'est un procédé coûteux mais avantageux à long terme !

Contrairement au procédé argentique, les images numériques peuvent être dupliquées INDEFINIMENT SANS PERTE

Sous forme numérique, les ouvrages peuvent être échangés, stockés, manipulés, traités et recopiés avec des coûts très réduits.

Les images numériques peuvent être reproduites sur de nombreux types d'affichage : papier (Book-on-demand), écrans (PDA, eBooks, poste PC...)

La numérisation est un projet **INTERDISCIPLINAIRE**

Nécessite des expertises (et donc des experts) sur

- Le choix des collections
- Les usages
- Le choix des métadonnées
- Les choix administratifs (accès par le web, poste de consultation sur place...)
- Les technologies informatiques (Matériel d'acquisition, logiciels d'analyse des images, formats de fichiers, les bases de données, l'interface du poste client, le langage de requêtes...)

Qui déterminent les fonctionnalités de la bibliothèque digitale

Quelles images pour quels usages ?

Pour être interprétées par une machine

Prévoir l'évolution des logiciels de reconnaissance

Qualité des images pour la reconnaissance *comment éviter les pertes d'information ? Comment réduire l'accumulation des bruits ?*

Pour être lues sur un écran

Confort visuels *quels type d'affichage ? Quels seront nos comportements demain (ex: e-book) ? Quels traitements pour améliorer la qualité des textes sur écran ?*

ATTENTION : une image bien lisible par l'homme n'est pas nécessairement bien interprétable par une machine et vice et versa !

Quelles images pour quels usages ?

Pour être lues sur le Web

- Consultation à distance* *Quelle compression adaptée au débit du réseau ?*
- Quelle interface de navigation ? comment éviter la copie ?*

Pour le stockage

- Conservation sur des supports compacts* *quels supports ? pérennité des supports et lecteurs ?*
- Problème du volume des données* *quelle compression ?*
- Qualité des images pour des traitements ultérieurs* *Quelle résolution ? Nombre de couleurs ? Quels traitements sont réversibles ?*

Pour la reproduction (Sous d'autres formats)

- Nécessite une totale rétroconversion*

Quelles images pour quels usages ?

Pour l'impression (livres rares à la demande, réimpression de titres épuisés, impression à domicile d'ouvrages achetés sur Internet)

📖 *Copyrights ? Marquage des images ?*

📖 *Résolution minimale requise pour l'impression ?*

☐ *Quelle restauration pour améliorer la qualité d'impression ?*

Pour l'enrichissement (Outils de recherche, navigation transversale par thème, auteur, date..., mise en commun des critiques, annotations, des différentes interprétations ou traductions d'un même ouvrage)

☐ *Rétro-conversion* ☐ *Quel niveau de rétro-conversion est possible pour un document donné ? Quel niveau d'automatisation ?*

☐ *Indexation* ☐ *Méta-Données ? Création automatique de liens ?*

☐ *Format de données* ☐ *Quel format est le plus adapté ?*

Numérisation et usages

Les choix techniques dépendent des usages que l'on fera de la bibliothèque numérique

Ces usages sont **multiples** suivant les époques et les sciences qui s'y rapportent :

Philologie (science de l'interprétation des documents)

Épigraphie (science des inscriptions)

Paléographe (science des écritures anciennes)

Codicologie (science des livres et de leurs supports)

Linguistique (grammaire, lexicologie ...)

Histoire des textes

...



**Comment Numériser ?
Comment prétraiter les
images ?**

**Acquisition, restauration
Traitements, compression
1ères étapes**

Qu'est ce qu'une image numérique ?



10	27	33	29
27	34	33	54
54	47	89	60
25	35	43	9

Une **image numérique** est une représentation **discrète** des **rayonnements visibles**.

Image = tableau[x,y] de pixels y=1..nb_lignes
x=1..nb_colonnes

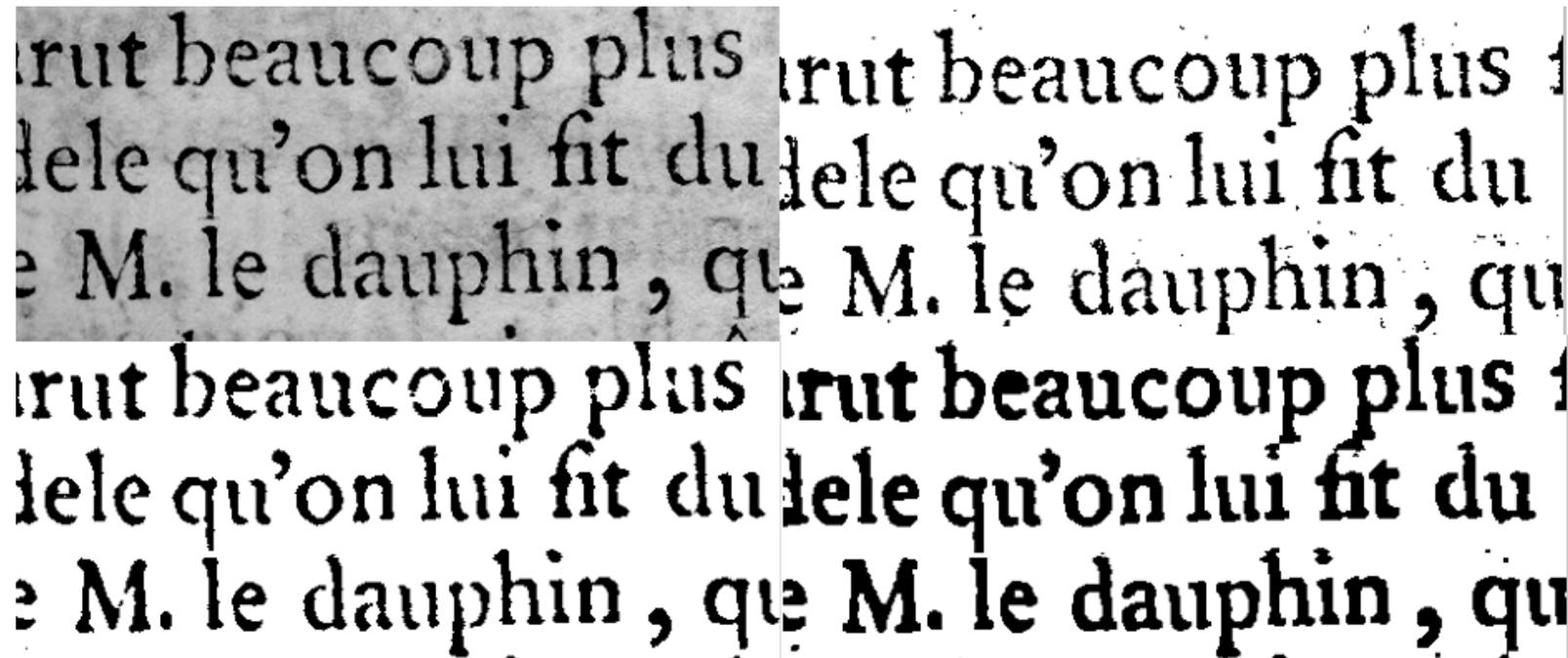
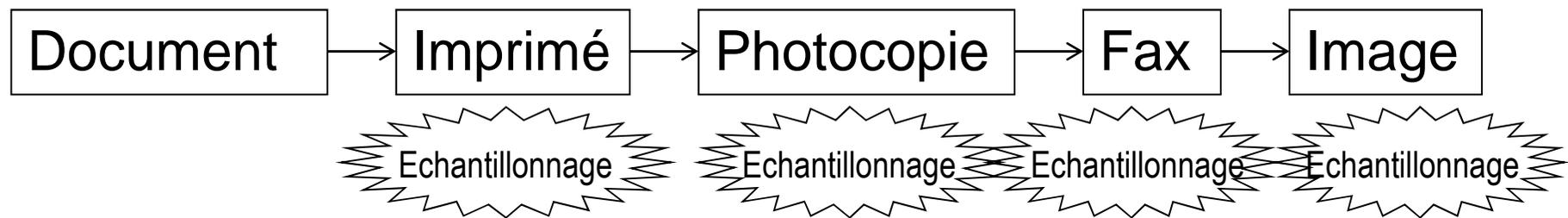
Chaque indice (X,Y) correspond à une information élémentaire de l'image nommée **pixel**.

La valeur de chaque pixel mesure le rayonnement quantifié sur cette surface élémentaire par un capteur. Celle-ci est quantifiée sur M valeurs.

Remarque : le concept de pixel est utilisé pour représenter indifféremment sa valeur ou sa position.

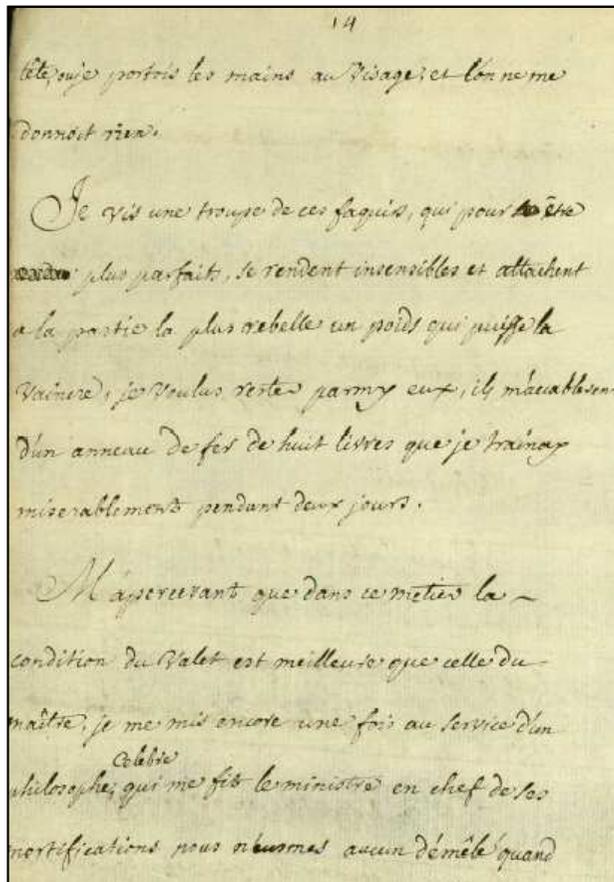
Numérisation des documents

Une image numérique de document est une représentation plus ou moins fidèle suivant la résolution du scanner, le nombre de couleurs utilisé et le nombre de transformations qui précèdent.

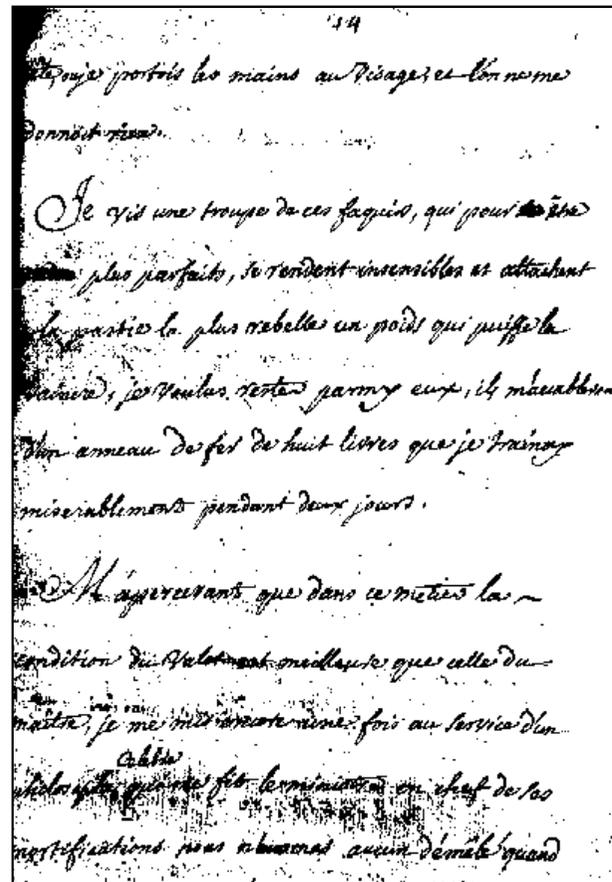


Numérisation des documents

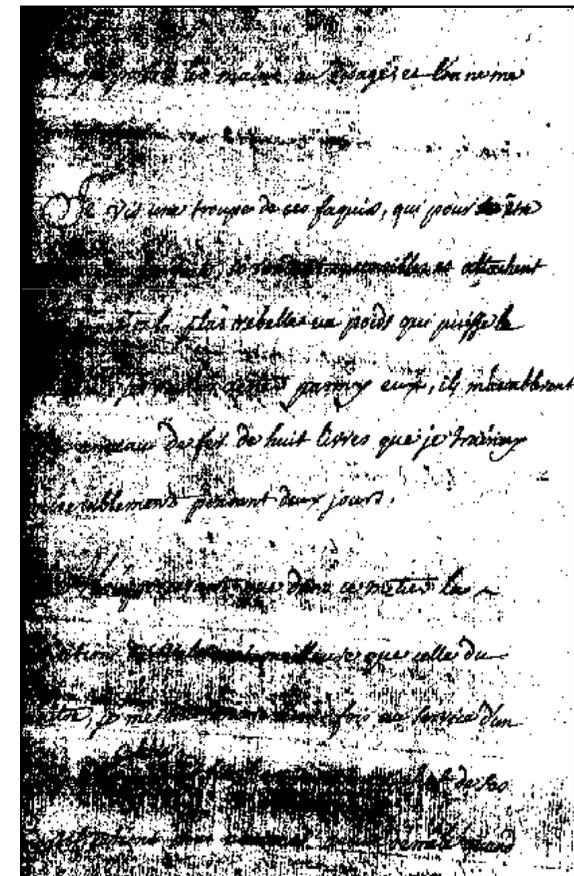
Image 24 bits



Photocopie



Fax

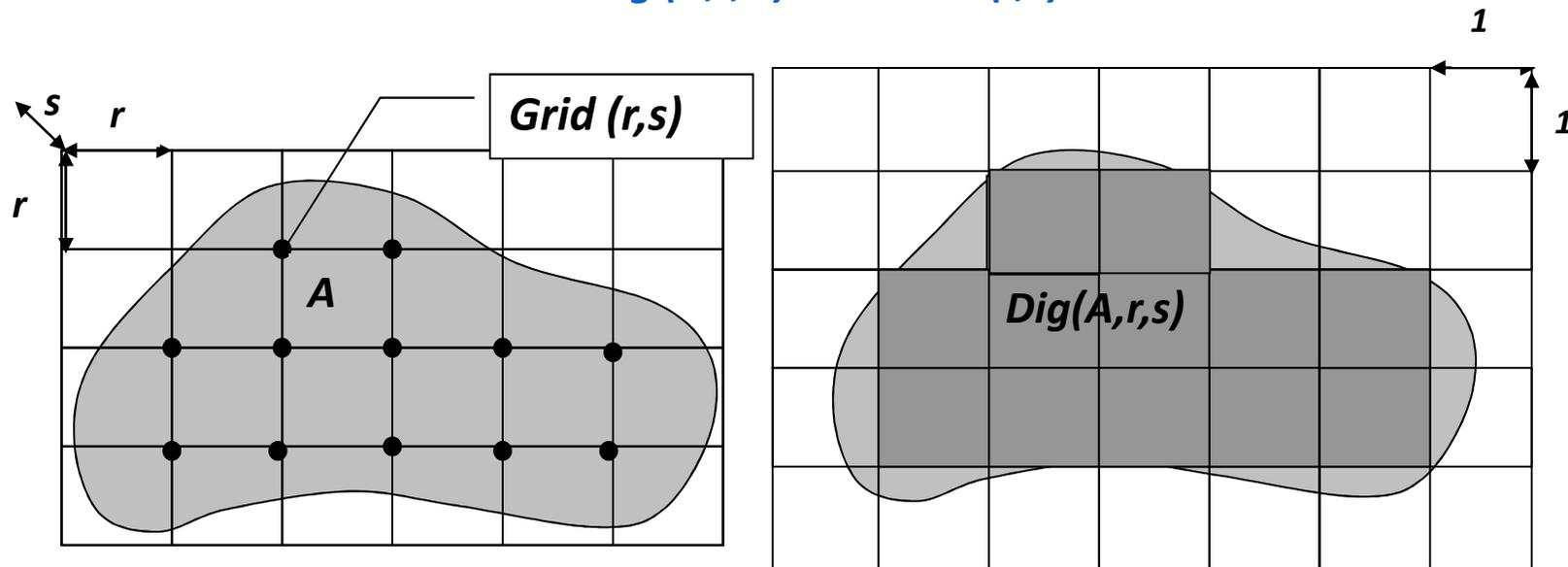


(Extrait collection Montesquieu, Histoire Véritable, 1750)

Géométrie discrète : du continu au discret

La **numérisation** est le résultat d'une discrétisation d'une forme continue **A** par une grille d'échantillonnage de taille r et qui se déplace avec un décalage s avec $|s| < r$

$$Dig(A,r,s) = A \cap Grid(r,s)$$

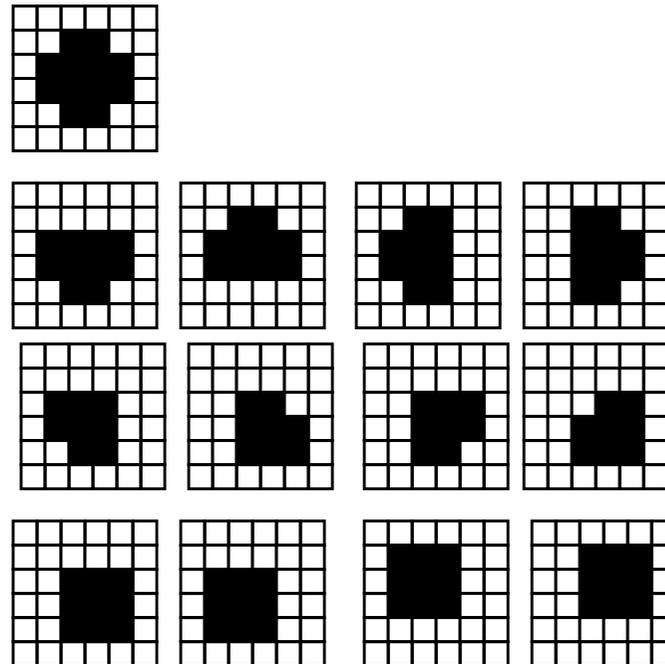
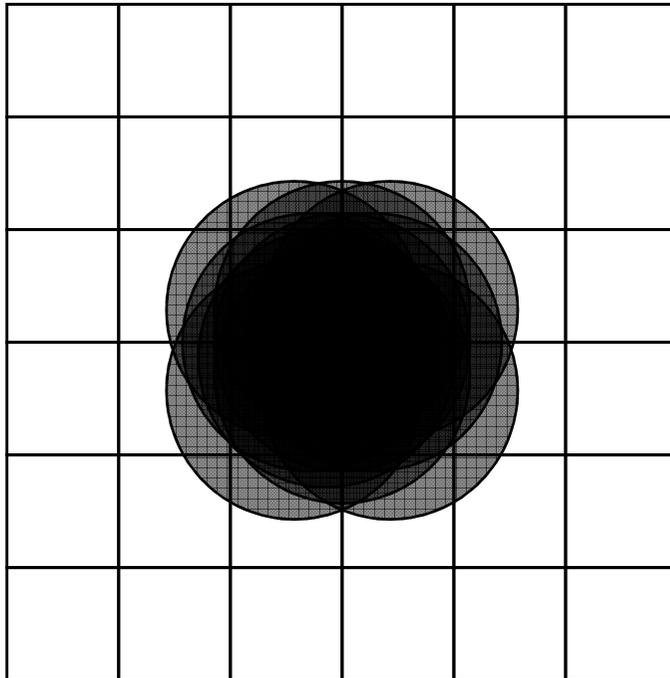


La grille d'échantillonnage correspond au capteur de la caméra

- ▶ La résolution r est la taille d'une cellule photosensible
- ▶ Le décalage s provient du placement imprévisible de la caméra

La théorie de la numérisation (1)

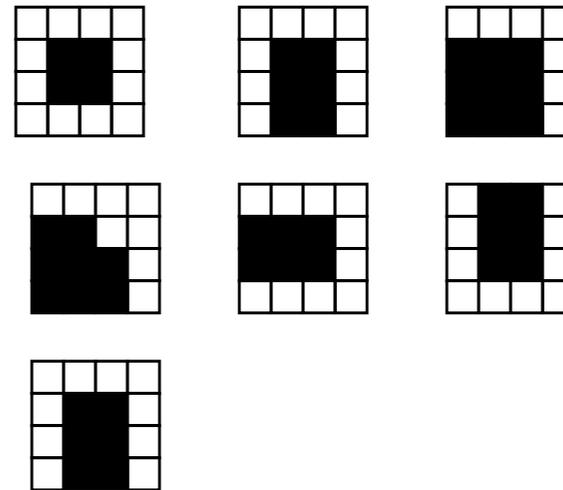
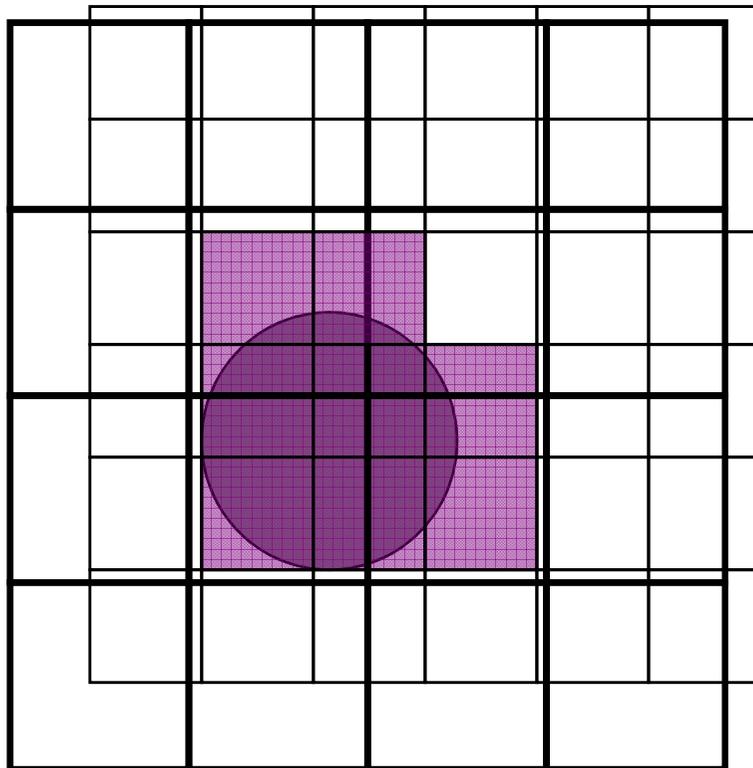
Combien de formes binaires obtient-on en déplaçant la grille de numérisation ?



La théorie de la numérisation (2)

Pour une forme numérisée avec la grille précédente, combien de nouvelles formes trouve-t-on avec une seconde grille de numérisation qui a un pas d'échantillonnage différent ?

(exemple : scanner un document imprimé ou photocopier un fax...)



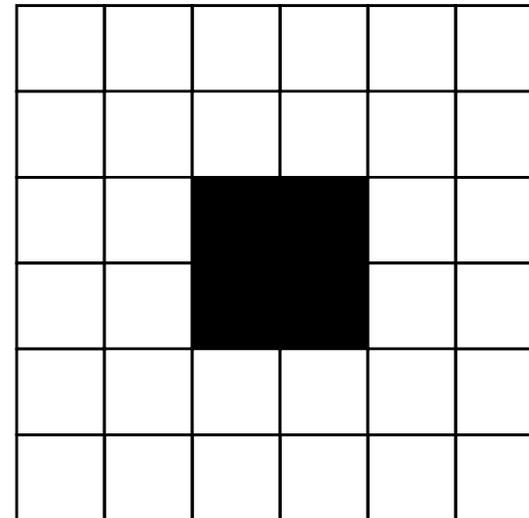
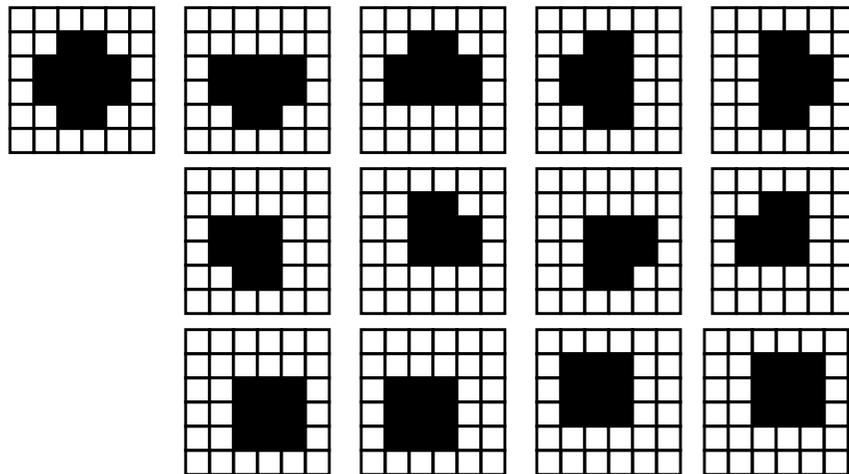
Pour chaque forme obtenue précédemment, on doit multiplier par le nombre de nouvelles formes trouvées avec la seconde grille de numérisation !

La théorie de la numérisation (3)

Le nombre de combinaisons de formes possibles augmente en fonction :

- de la complexité des formes,
- du pas d'échantillonnage de la grille
- du nombre de grilles superposées

Remarque : Certains pixels apparaissent toujours dans toutes les grilles !



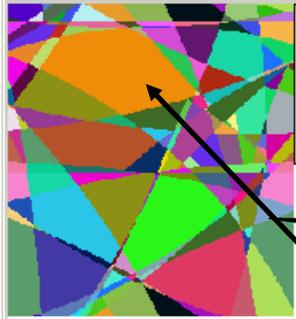
Pixels toujours présents

La théorie de la numérisation (4)

La numérisation ne conserve pas la topologie des formes :

- ▶ Le 'e' de police 'Times' dans une grille 10x10 produit
 - 196 formes différentes
 - 50% font apparaître une rupture de tracé ou un remplissage des boucles.
 - 20% ressemblent plus aux caractères 'c' et 'o' qu'au caractère 'e'
- ▶ Fréquences d'apparition de chaque forme ?

Le Modulogramme permet à la fois de connaître le nombre de formes et la fréquence d'apparition de chacune d'entre elles



Forme continue originale

Modulogrid

pattern=196 Nombre de tirages =38799

Nombres de formes différentes

Forme numérisée obtenues et sa fréquence d'apparition

	3351	306	183	119	7	62	168	241	670	619	115	189	82	19	32	25	138
	41	21	449	155	32	15	15	13	122	36	69	68	4	83	86	39	229
	97	6	18	248	400	3249	28	803	12	148	234	183	255	16	122	7	6
	311	163	10	86	7	70	155	29	267	56	1	248	62	1075	23	376	157
	77	27	392	47	578	2	54	432	5	62	6	2	43	63	3	128	4
	105	3	473	41	39	30	104	47	19	1627	10	64	49	314	85	97	43
	2	179	278	1	150	232	130	19	36	42	261	23	30	49	85	42	4
	295	2	134	11	126	1	38	141	9	88	37	4	1	990	314	181	1115
	2	79	157	141	79	14	28	2077	129	30	1476	15	39	215	157	256	23
	35	39	381	15	11	3	23	1384	1964	42	105	109	21	90	297	6	32

Images couleurs

Information Maximale (16 millions de valeurs possibles/pixels)



Pour un usage avancé et complet (codicologie, étude des miniatures, encres, décorations ; restauration&segmentation...)

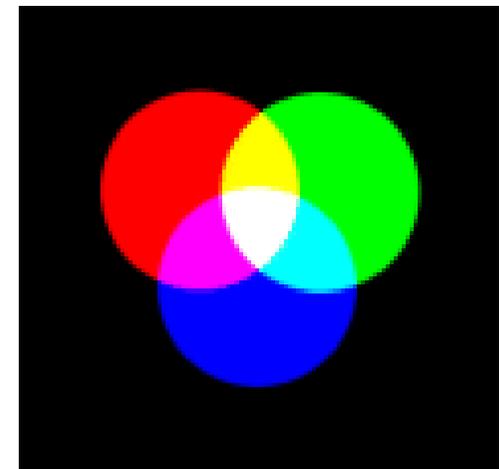
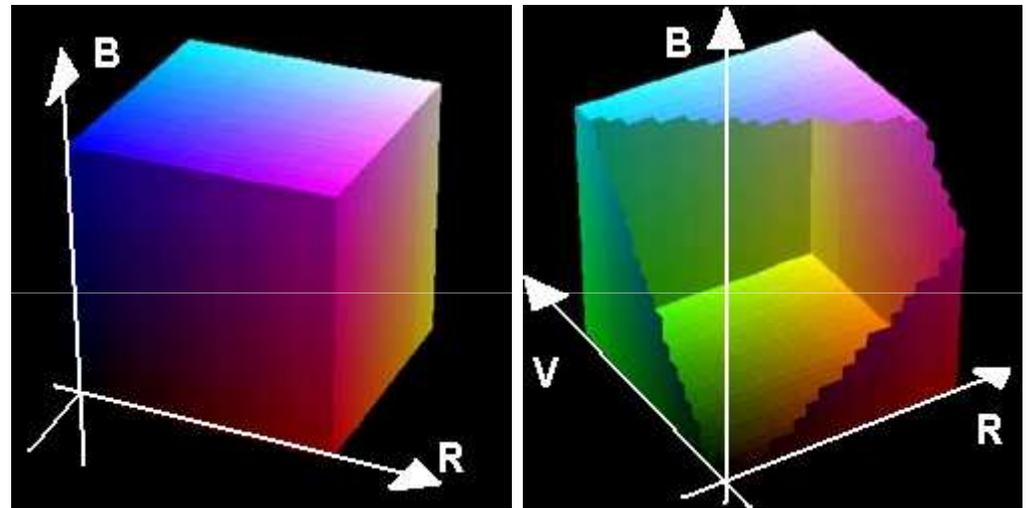
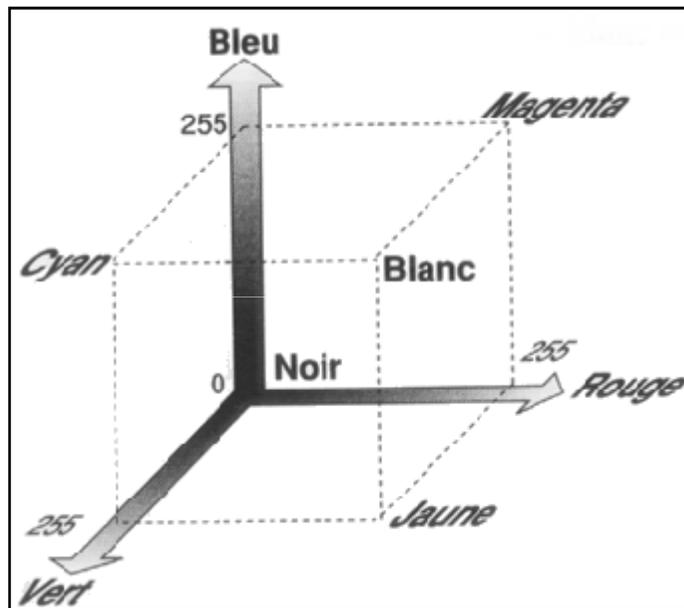
Images couleurs (2)

- L'œil humain ne distingue qu'une centaine de milliers de nuances de couleurs différentes
- On parle d'« intensité lumineuse » et de « chromaticité »
- Un pixel est codé par trois valeurs numériques
 - La signification de ces valeurs dépend du type de codage choisi
 - Le plus courant (stockage, photo et visu. écran) : RVB (ou RGB)
 - Quelques autres :
CMJ (ou CMY), TSL (ou HSL en vidéo), YUV, YIQ, Lab, XYZ ...
 - Exemple : Vraie Couleur RVB (True Color)
chaque valeur code une composante couleur sur 8 bits

donc $2^8 \times 2^8 \times 2^8 = 16\,777\,216$ couleurs possibles

Images couleurs (3)

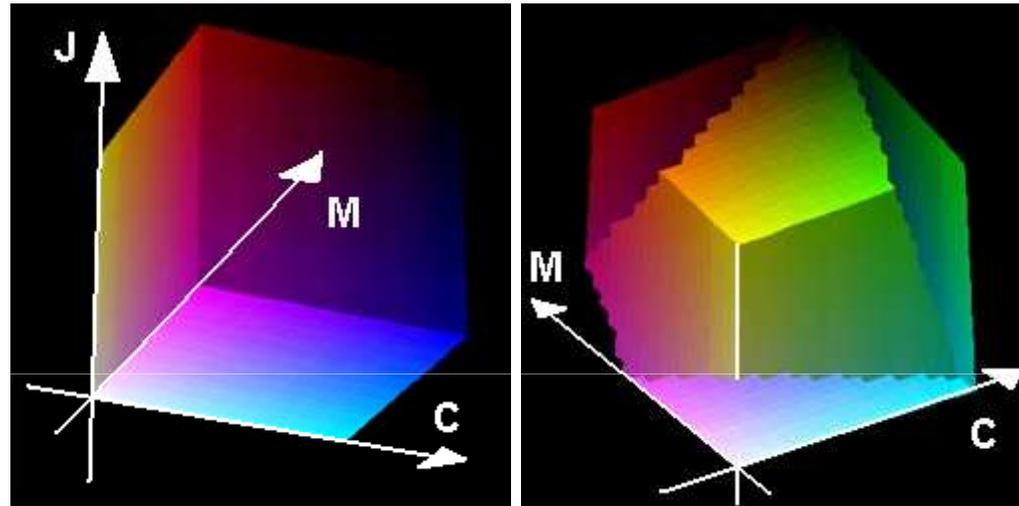
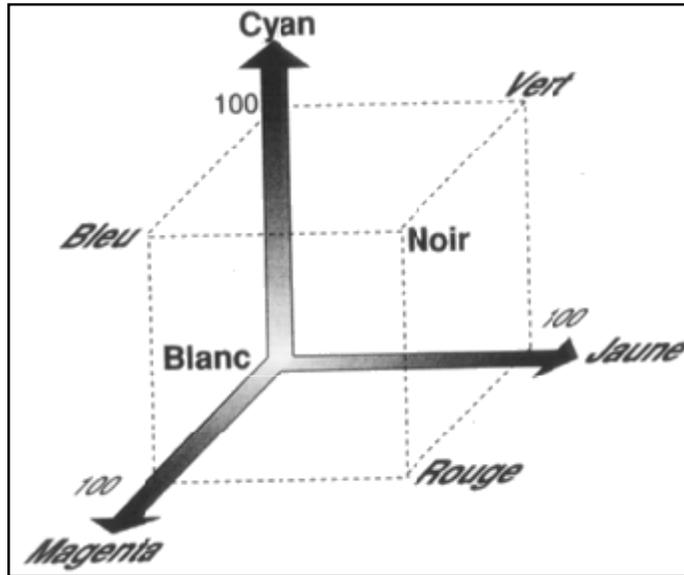
- L'espace couleur **RVB** (Rouge Vert Bleu)



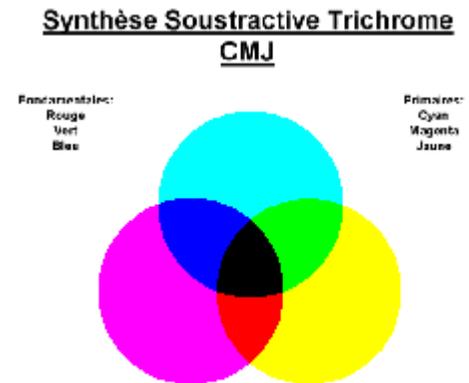
- Les couleurs sont contenues dans un cube
- La restitution des images sur écran utilise cette représentation RVB (synthèse additive)

Images couleurs (4)

- L'espace couleur **CMJ** (Cyan Magenta Jaune)



- Les couleurs sont contenues dans un cube
- Principe utilisé en imprimerie (par mélange d'encre)
- La restitution des images sur papier utilise cette représentation CMJ (synthèse soustractive)

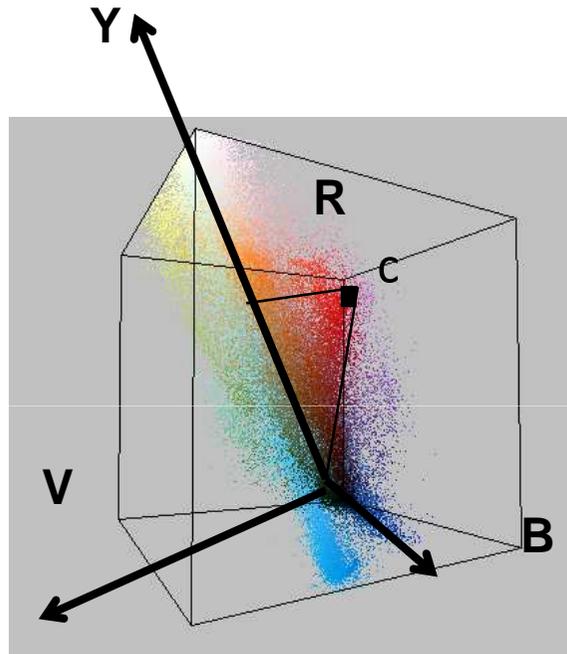
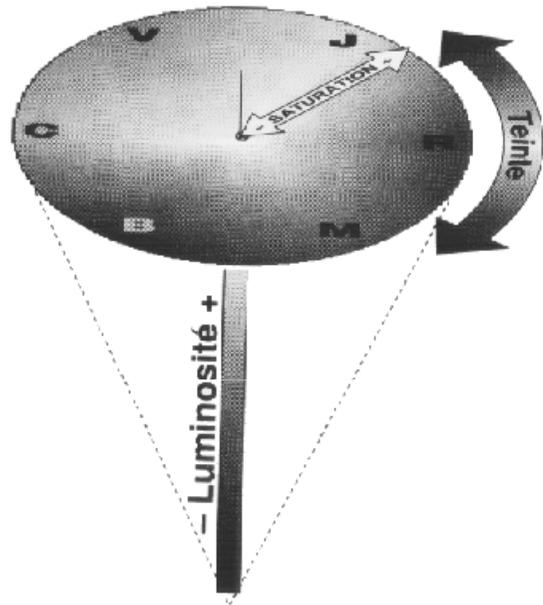


Images couleurs (5)

espace couleur TSL (Teinte Saturation Luminance)

$$c=(R,G,B)$$

Y= l'axe luminance



« La couleur
des peintres
»

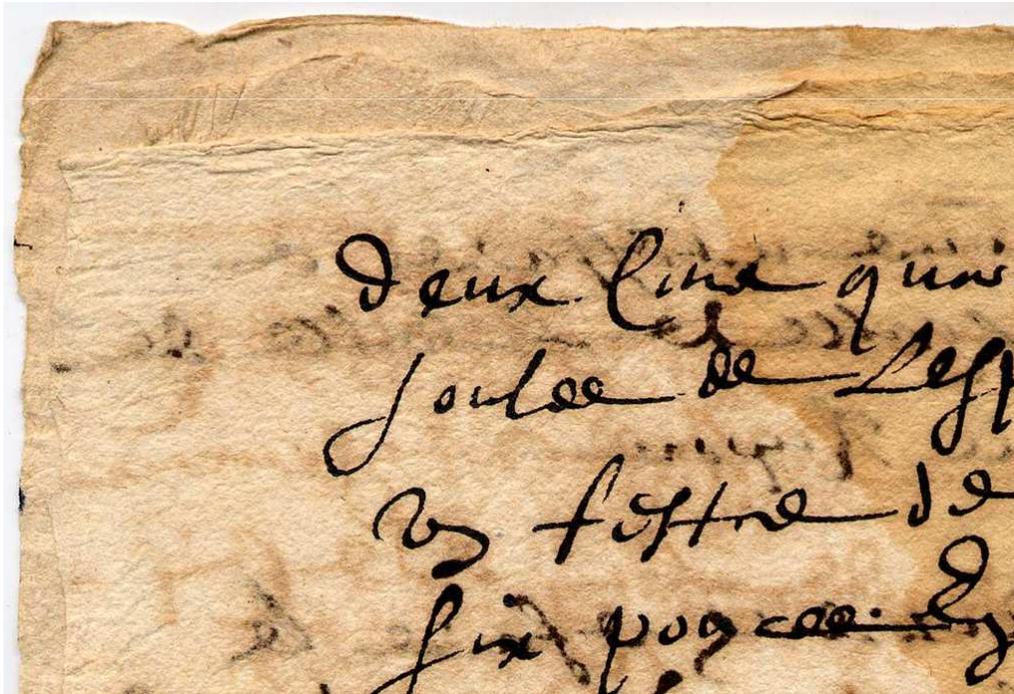
- ▶ **Luminance** = $\text{projection}(c/Y)$ « quantité de noir »
- ▶ **Teinte** = $\text{angle}(c,Y)$ « longueur d'onde de la couleur »
- ▶ **Saturation** = $\text{distance}(c,Y)$ « quantité de blanc », une couleur à 100% est dite pure, ne contient pas de blanc

Numérisation directe en niveaux de gris ou en couleur

Numérisation directe toujours meilleure que celle des microfilms

L'information niveaux de gris ou couleur est utile pour la segmentation

Suivant les usages , l'information couleur est indispensable



Exemple de la qualité d'image par une numérisation directe

La numérisation en Niveaux de gris (1)

Toujours plus d'information dans une image faible résolution en niveaux de gris que dans une image haute résolution en binaire

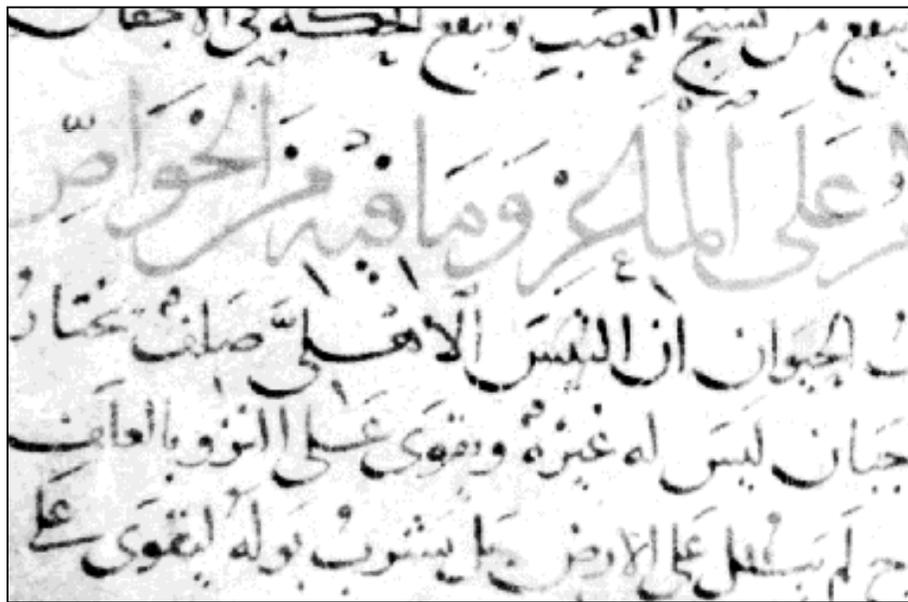


Image 16 niveaux de gris

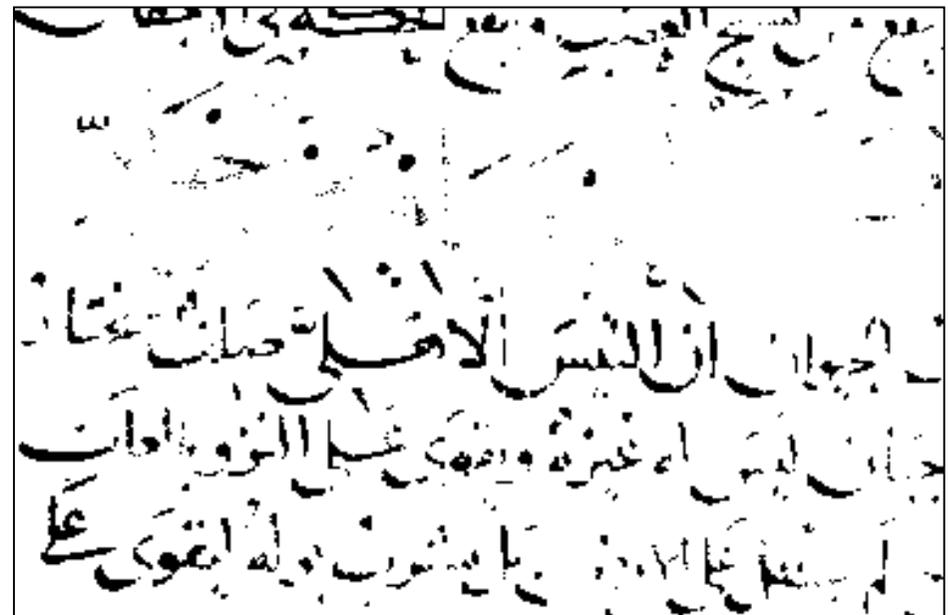


Image binaire

Images niveaux de gris (2)

information moyenne (256 valeurs possibles / pixel)



Pour un usage multiple (codicologie, étude des textes, restauration et segmentation limités des images ...)

Images niveaux de gris (3)

Représentation 3D d'une image à niveaux de gris

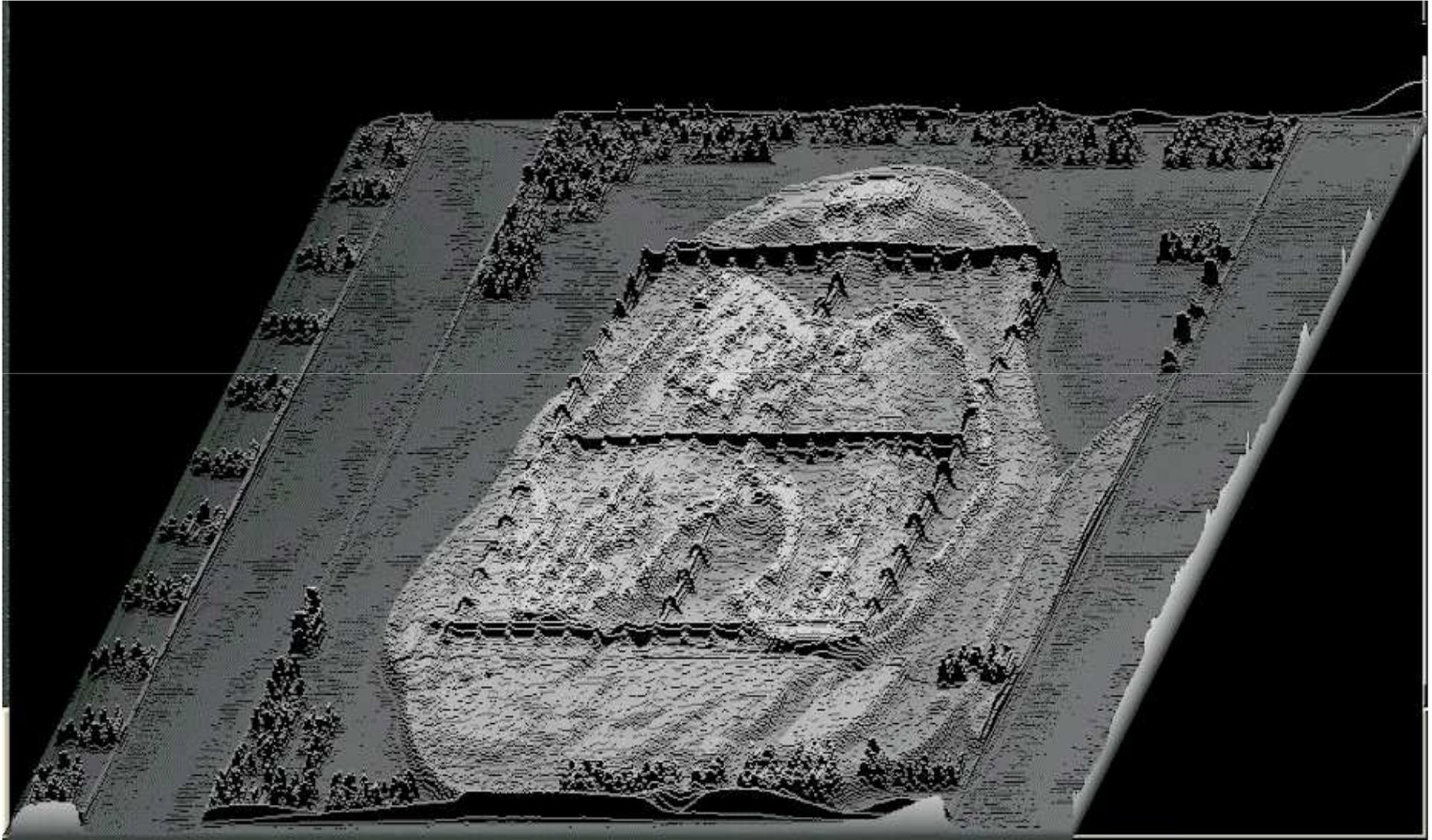
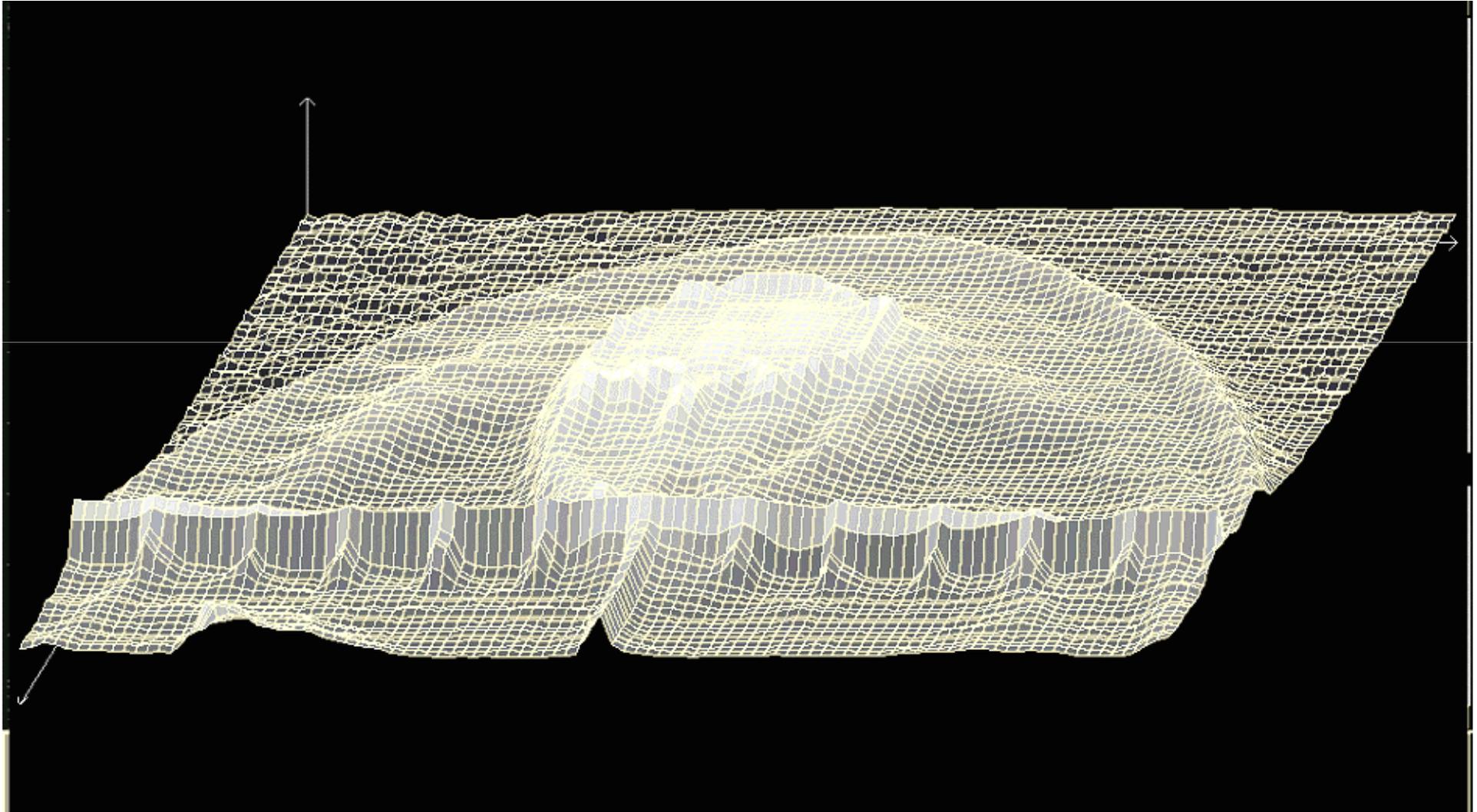


Image numérique vue représentée par une surface

Images niveaux de gris (4)

Mais c'est une surface discrète



Gros plan sur la région supérieure de l'image

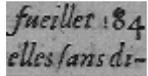
Images binaires (N&B)

information limitée : 2 valeurs possibles: 0 (blanc) 1 (noir)



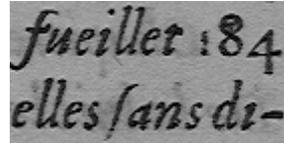
Binaire pour un usage limité à celui des textes (dans les limites de la qualité de la méthode de binarisation)

Choix de la résolution



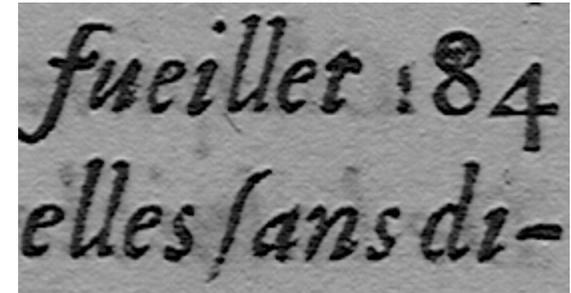
feuille : 84
elles sans di-

150 dpi



feuille : 84
elles sans di-

300 dpi



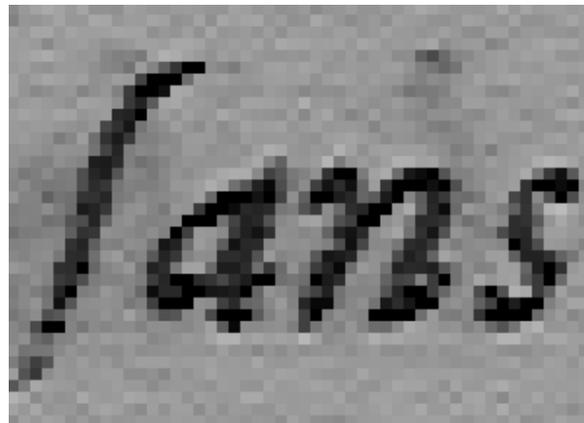
feuille : 84
elles sans di-

600dpi

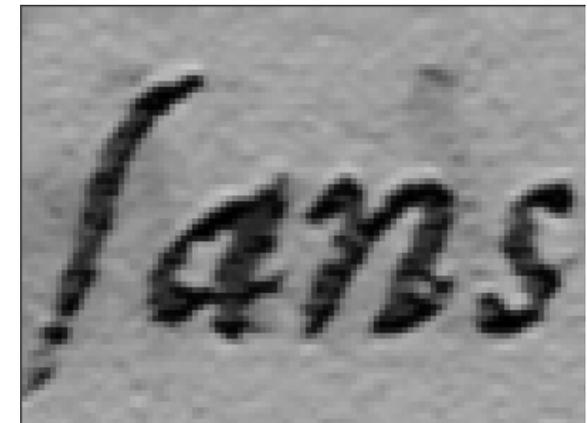
...Après harmonisation des tailles à l'écran par zoom :



sans



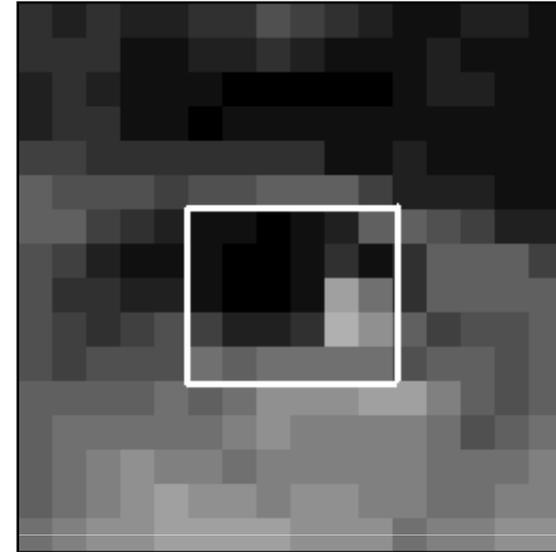
sans



sans

DPI = Dot Per Inch (nombre de pixels réels par pouce)

Choix de la résolution



Exemple :

Résolution : 180x180 pixels
avec 256 niveaux de gris
(16 réellement utilisés ici)

16	16	0	16	32	96
16	0	0	16	48	16
16	0	0	16	160	112
64	32	32	48	176	144
112	96	112	112	112	112

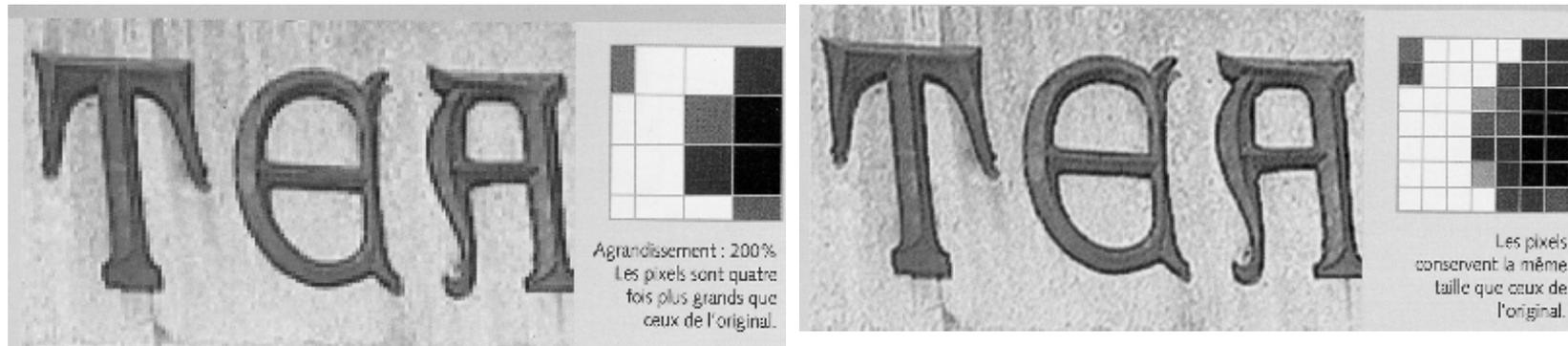
Règles de conduite (1)

- ▶ Toujours utiliser la résolution la plus élevée possible!
- ▶ C'est à dire plus d'information...

Tous les systèmes d'analyse de documents numérisés (*reconnaissance des caractères, des structures, indexation...*) sont sensibles à la résolution des images

La vision humaine n'est pas sensible à la résolution des images !

Exemple :



Résolution : Donné en dpi ou ppi (nb de points par pouce) il définit le nombre de pixels réels par unité de mesure

Règle de conduite (2)

- ▶ Numériser avec la mire couleur et une échelle
- ▶ Equilibre entre la qualité des images et la vitesse de capture

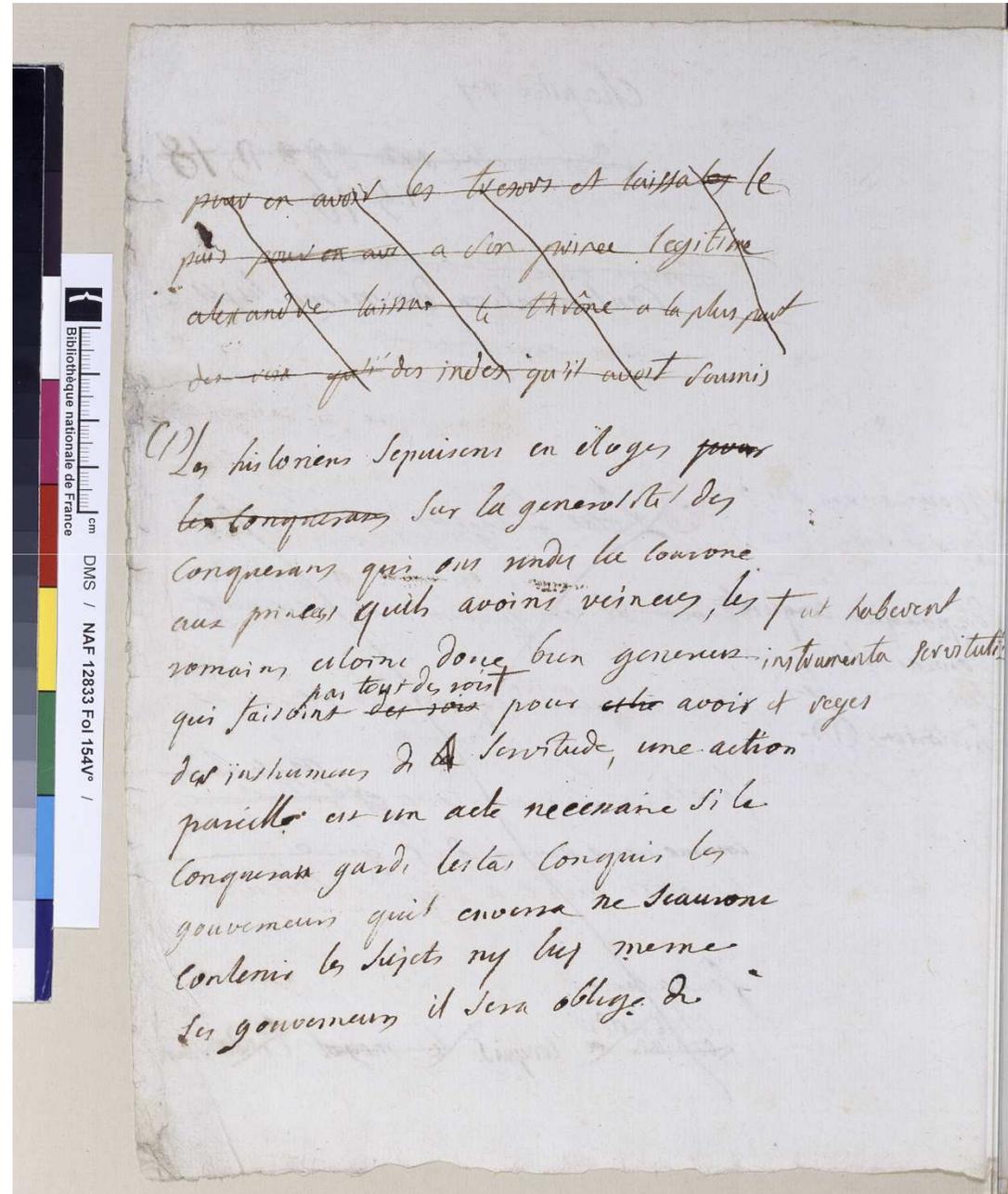


Bib. Mazarine

Règle de conduite (2b)

- Numériser avec la mire couleur et une échelle

Feuillets Montesquieu
BM Bordeaux



Règle de conduite (3)

- ◆ Fixer le facteur d'échelle inférieur ou égal à 100% (éviter les agrandissements lors de la capture)
- ◆ Ne jamais sur - échantillonner une image !!
- ◆ Seule la résolution optique d'un scanner est importante (pas la résolution interpolée)

Règles de conduite (4)

► Eviter les étapes intermédiaires

Mauvais exemples:

document papier → microfilm → scanner

document papier → photographie → scanner

document électronique → imprimerie → scanner

document papier → photocopie → scanner

Bons exemples :

document papier → scanner

document papier ou objet → photographie numérique

signer une convention pour récupérer les documents électroniques chez les éditeurs

Quelle qualité d'images pour quels usages

Avant de commencer la numérisation, il faut faire une étude préalable

- Numérisation en interne ou sous-traitance ?
- Pour quelle utilisation numérise-t-on ces fonds ?
- Quelle résolution spatiale ?
- Numérisation en binaire ? Niveau de gris ? Couleur ?
- Quelle échelle ?
- Quels traitements pour restaurer les images ?
- Quel format de fichier (avec/sans compression avec/sans perte)
- Quel support pour les sauvegardes (durée du support dans le temps)

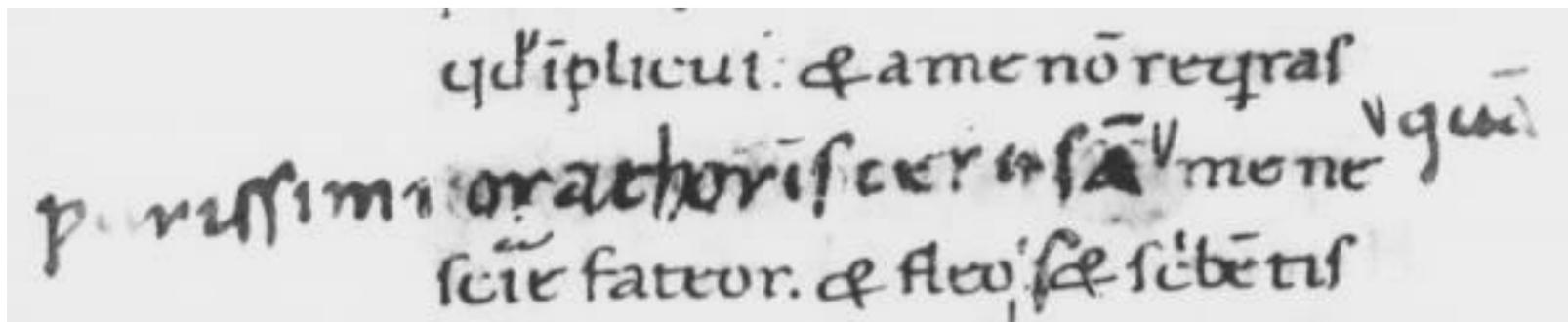
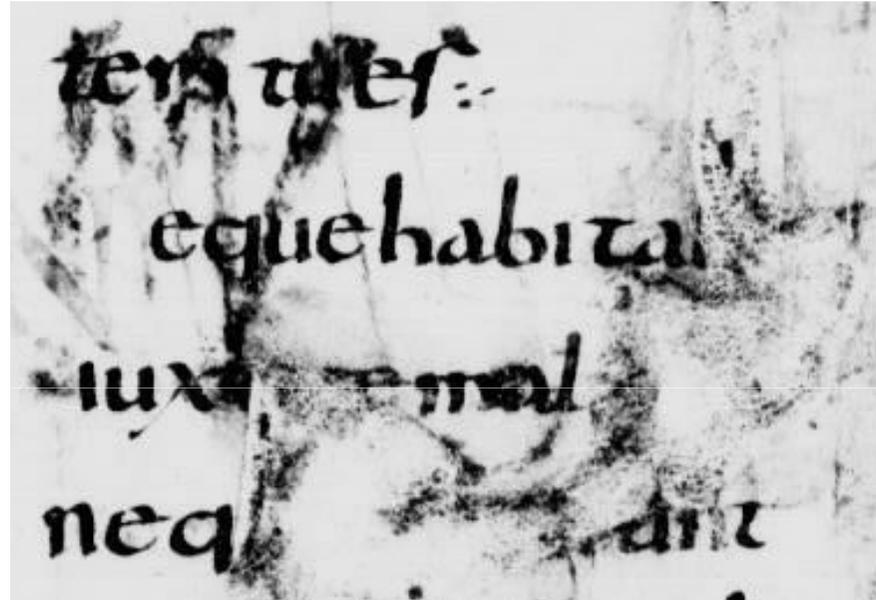
Critères de choix du matériel

- 1) **Préservation des originaux** (*Scanners sans vitre, Lumière froide, plateau pour préserver la reliure du livre..*)
- 2) **Qualité des images** (*Résolution et fidélité de restitution des couleurs*)
- 3) **Vitesse** (*influence les coûts et le temps*)
- 4) **Automatisation, service et logiciels** (*Ergonomie zero-fatigue*)
- 5) **Coûts** (*amortissement , prix par page...*)



La qualité des images

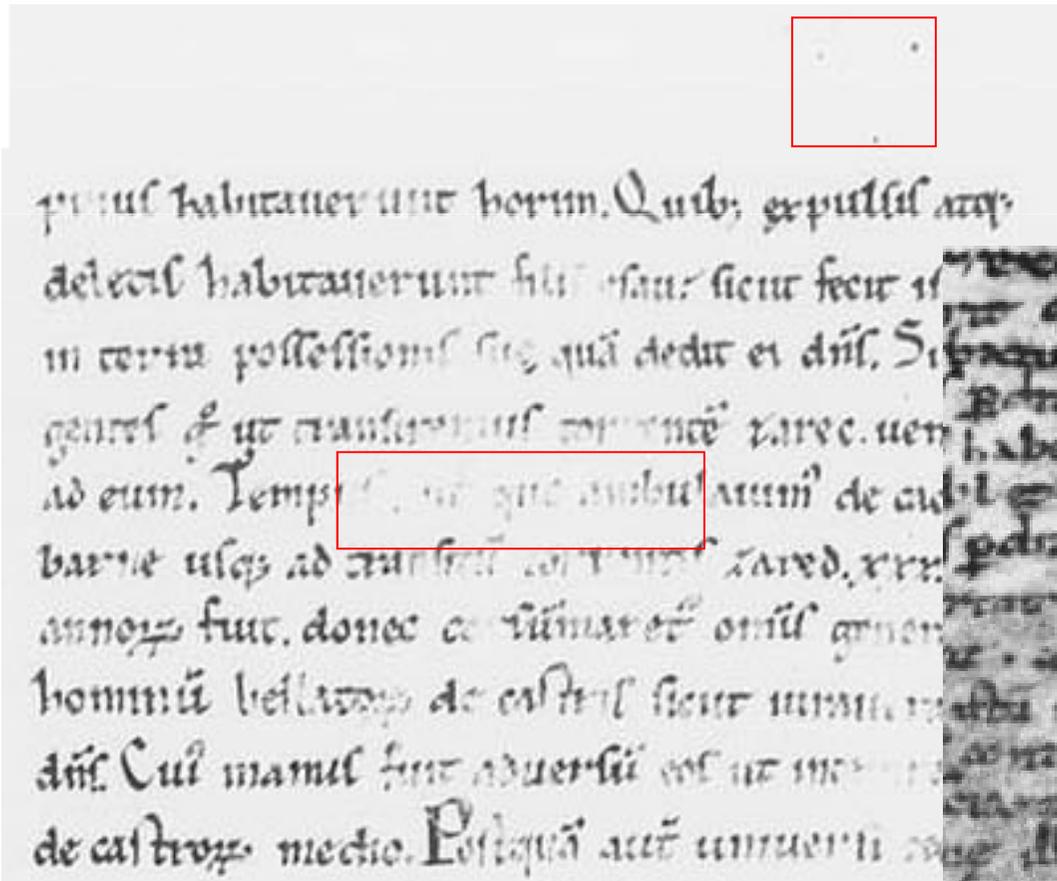
Documents dégradés avant la numérisation



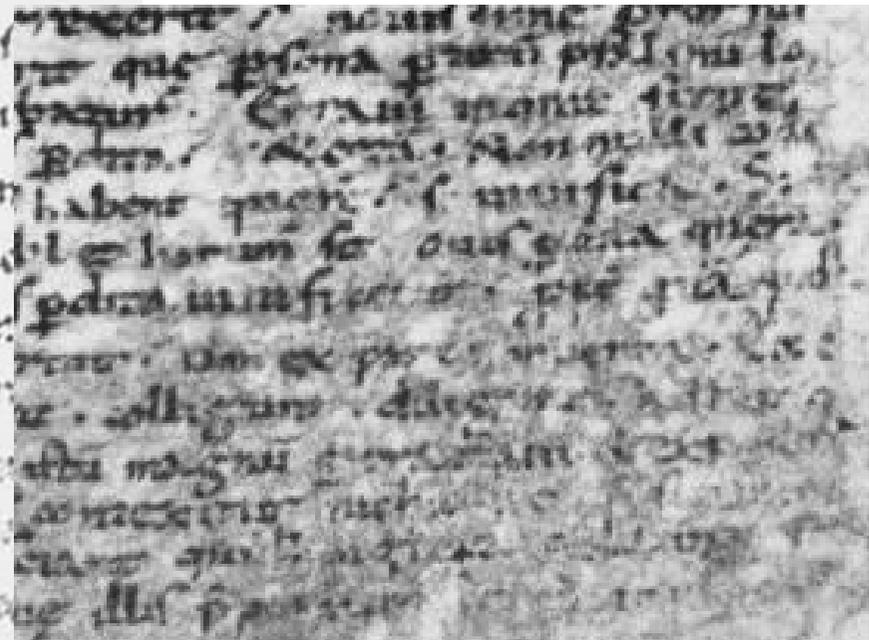
Numérisation de microfilms

Plus rapide et plus économique qu'une numérisation directe
...mais aux détriments de la qualité des images

Le microfilm nécessite le rehaussement du contraste qui efface les
niveaux de gris intermédiaires



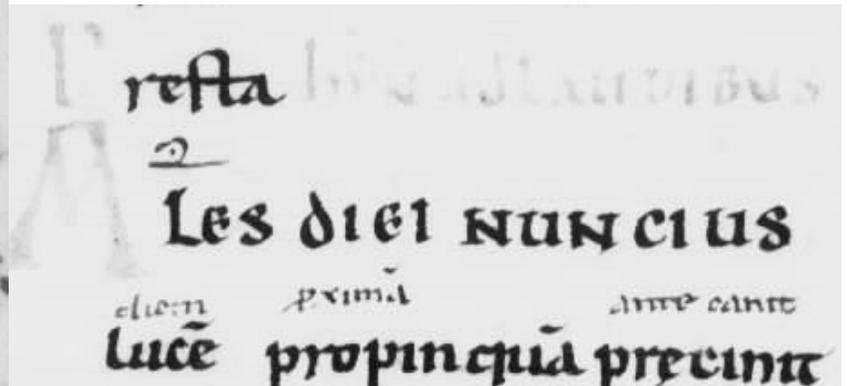
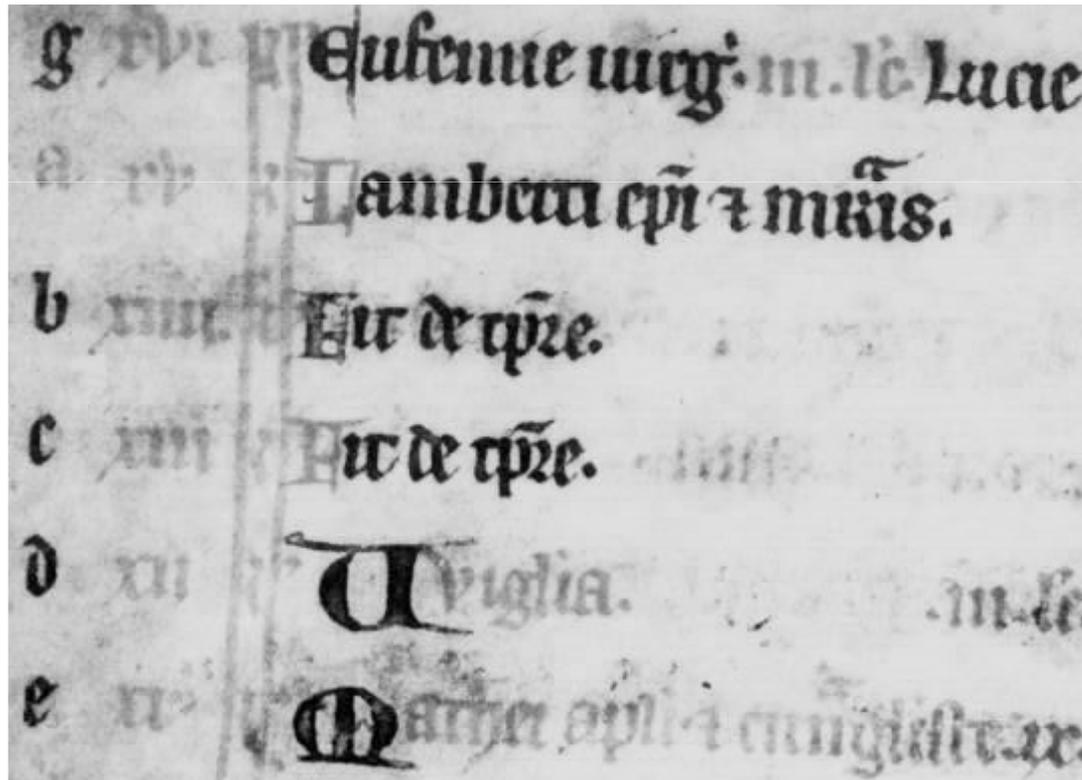
Caractères effacés



Numérisation de microfilms

Perte d'information sur la couleur des caractères

Le processus de microfilmage estompe les caractères rouges
AVANT même la numérisation

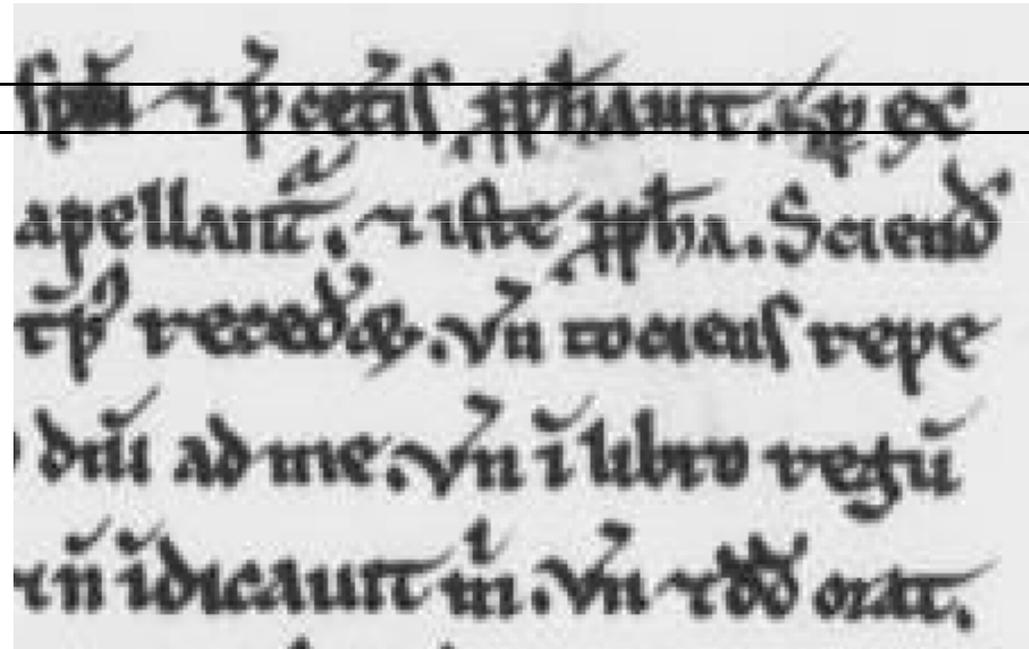
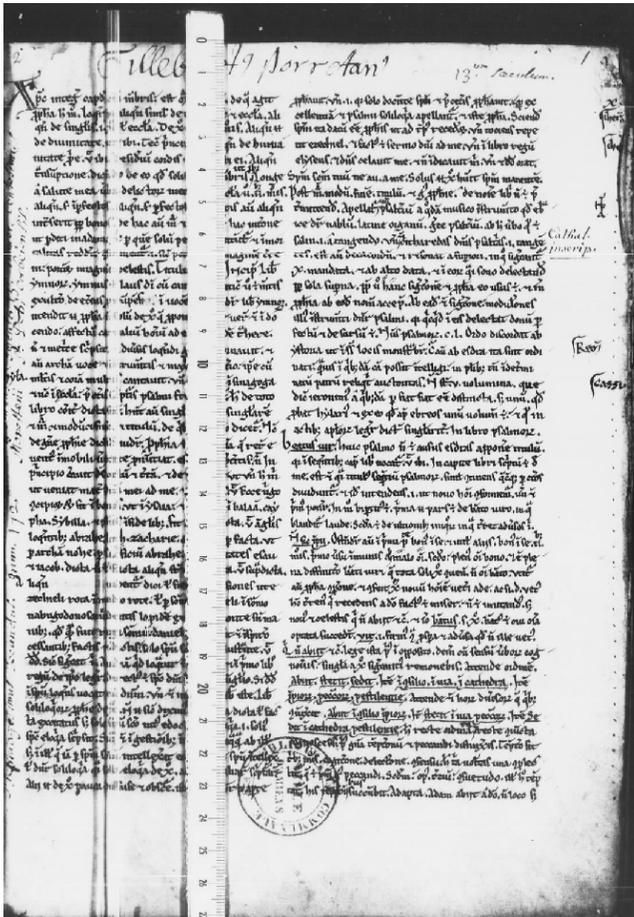


Numérisation et résolution

Choix de la résolution dépend de la taille (Hauteur) du plus petit caractère présent dans le document.

H > 12 pixels pour la lisibilité

H > 24 pixels pour une machine !



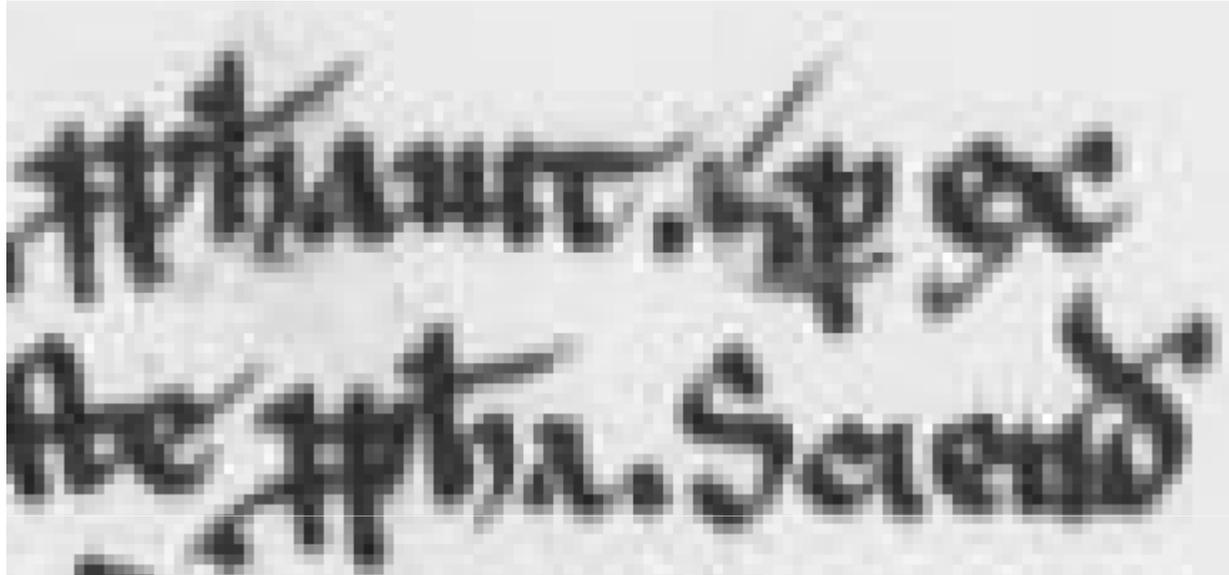
Résolution efficace

=

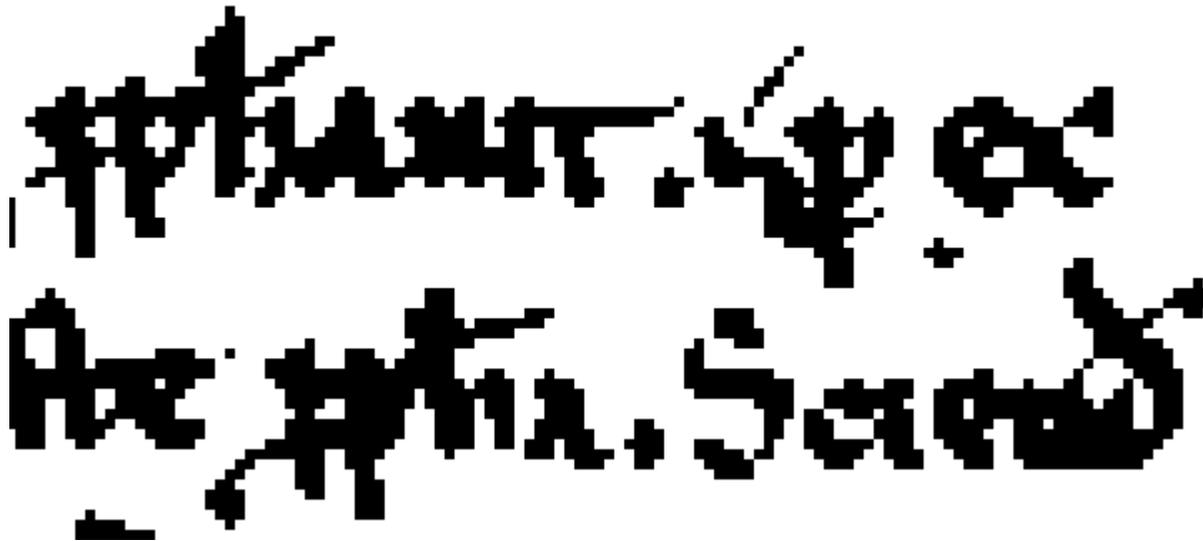
Nbr de pixels réels / Hauteur en pouces

En basse résolution

Il est impératif de numériser en niveaux de gris !



Résolution insuffisante
mais reste lisible à
l'Homme grâce aux
niveaux de gris



En binaire le texte n'est
plus lisible par l'Homme et
par la machine

**Le problème c'est que
la machine binarise
l'image pour l'analyser**

Problèmes de résolutions

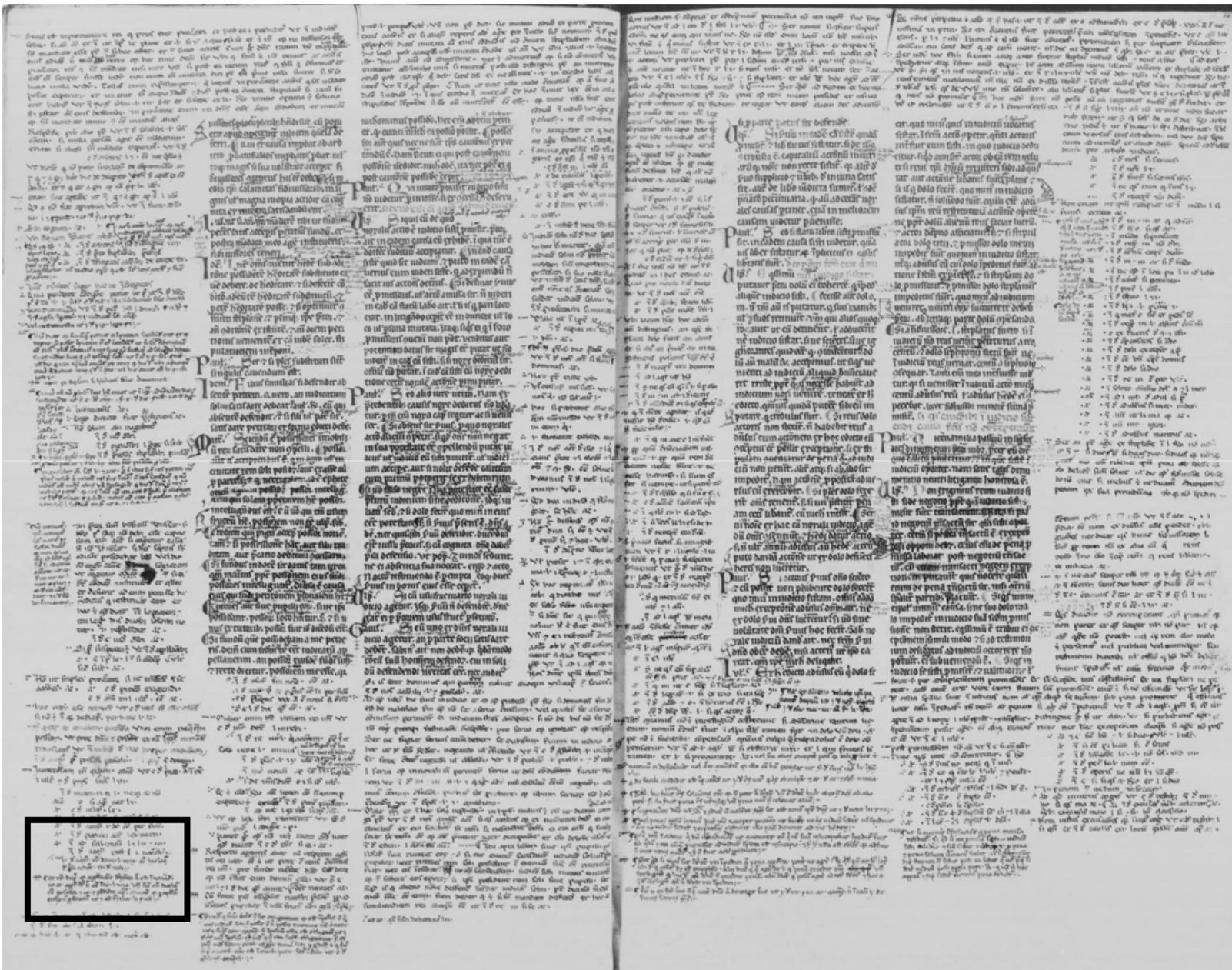
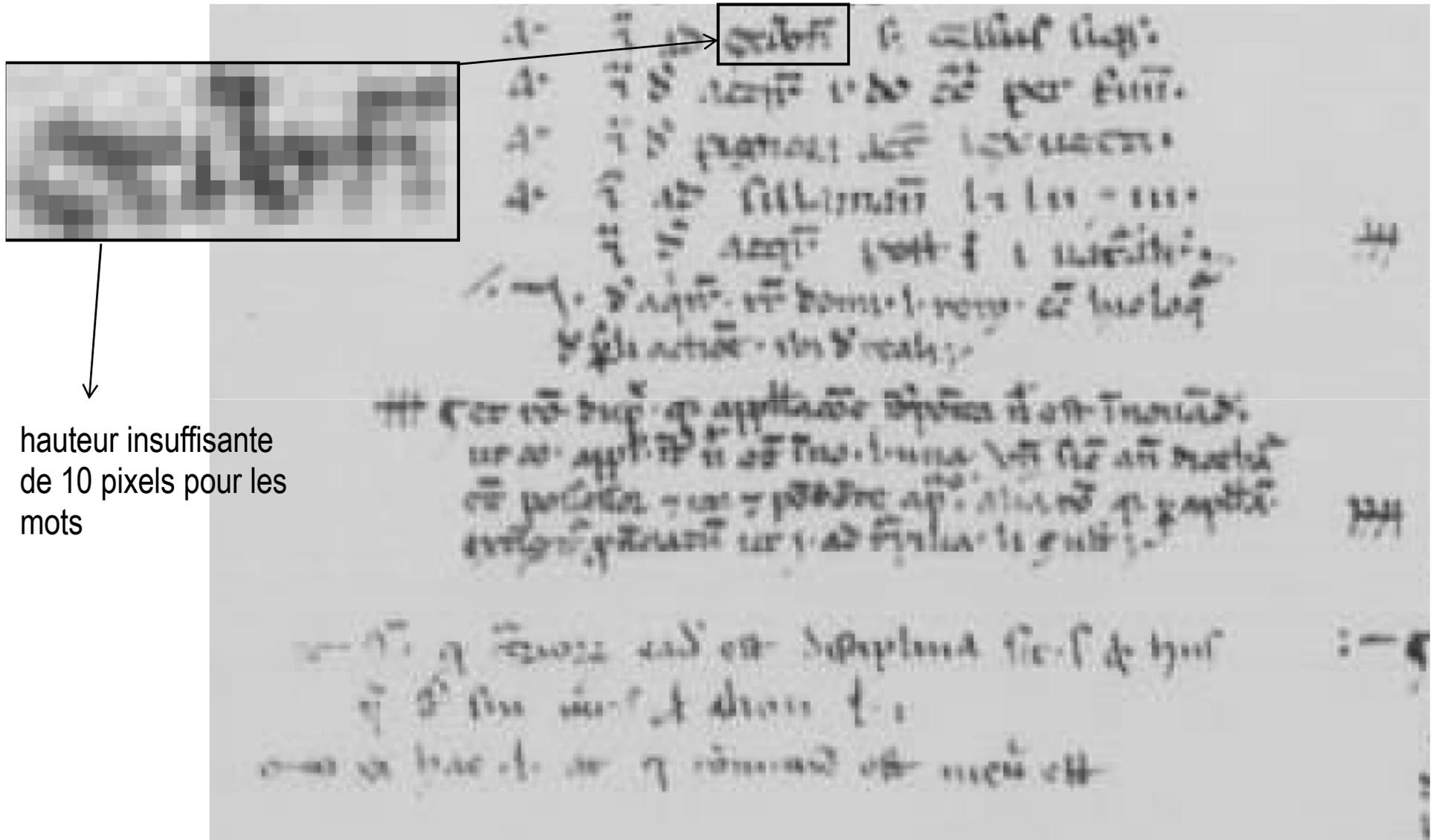


Image 2000x1600 issue de microfilm et compressé JPEG 60%

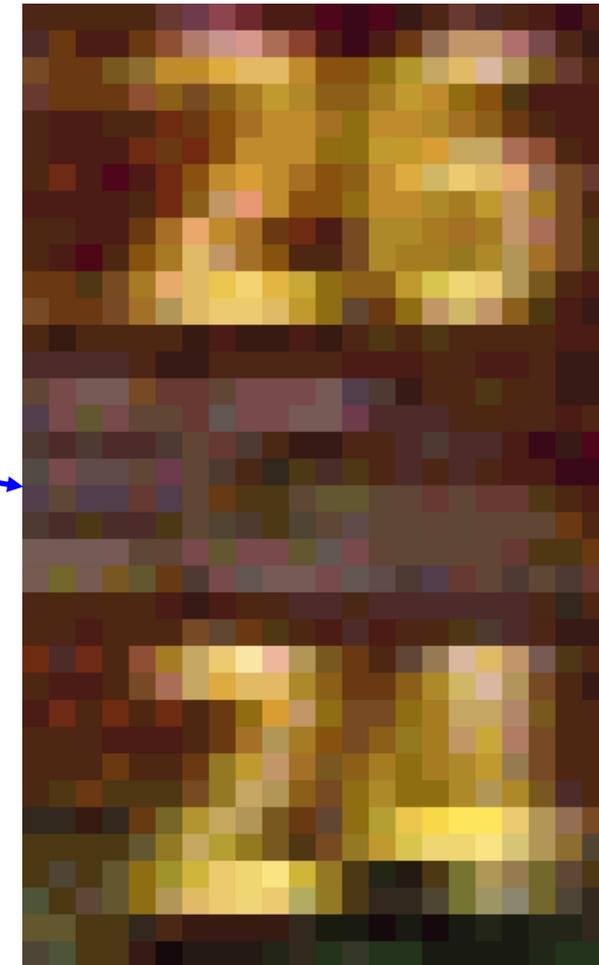
Problèmes de résolutions

Impossible à lire même pour un expert alors pour une machine...



Zoom de la portion d'image 330x220 du coin en bas à gauche

Exemple en vidéo ...



Que faire avec une image numérisée?

.Peut on y accéder?

- Recherche
- Navigation
- Lecture

.Peut-on l'indexer et la récupérer sur simple demande?

Dans sa représentation initiale, pas réellement!

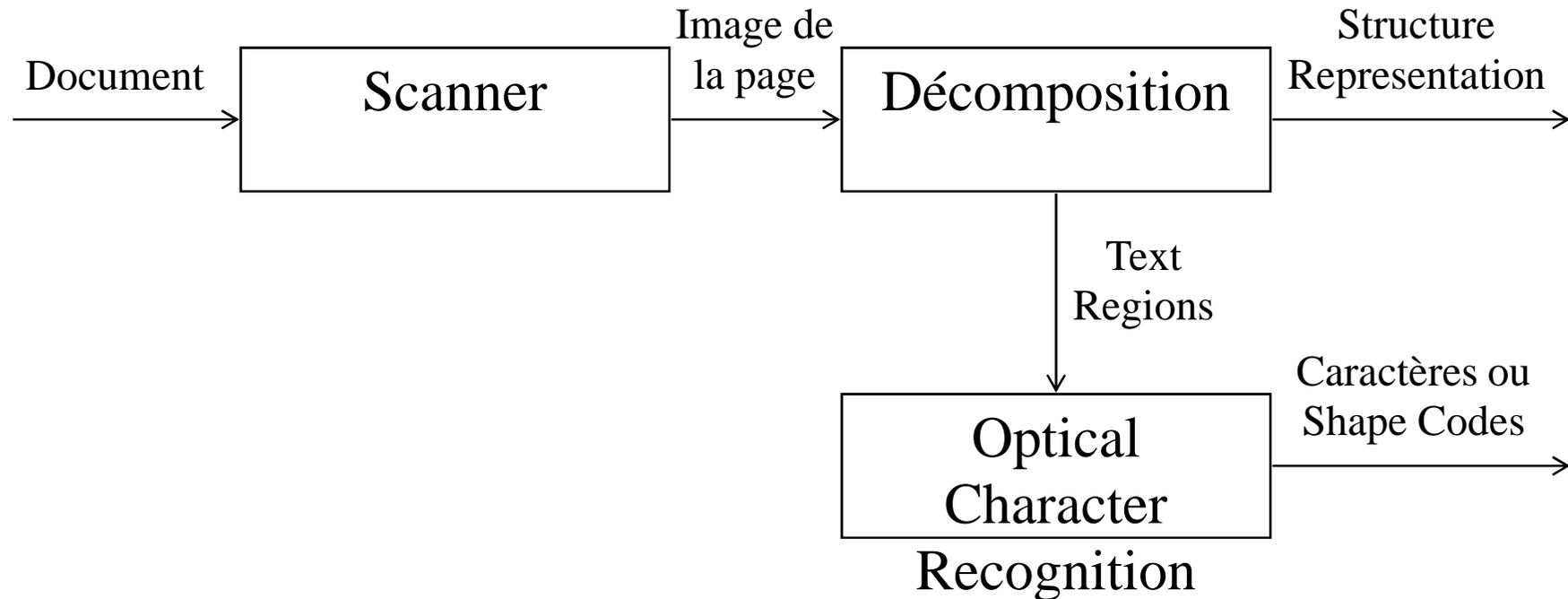
.Il est néanmoins possible de

- La visualiser
- L'imprimer
- La recopier... mais pas beaucoup plus

Quels outils pour quels usages ?

Comment accéder aux textes des documents numérisés ?

Principe général



Comment accéder aux textes des documents numérisés ?



Acquisition
Scanner / Caméra

Amélioration de
l'image

Imagerie

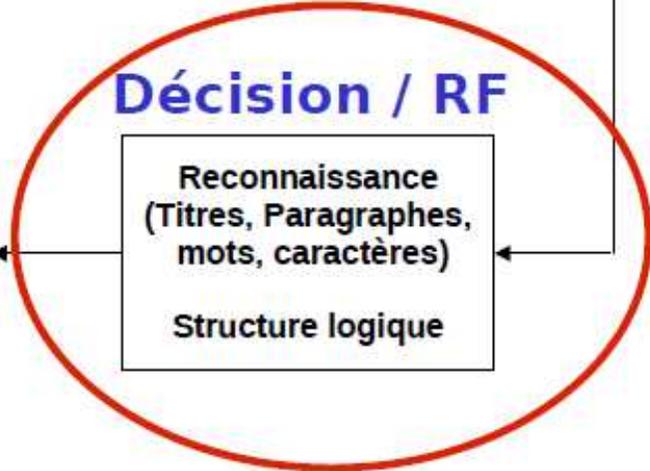
Localisation
(Blocs de texte,
ligne, mot)

Structure physique

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<!-- GARDAVE_ML.DTD dedicated to the description of manuscripts layout -->
<ELEMENT transcription (Bloc)*>
<!-- a transcription is made up of one "bloc" entry at least -->
<ELEMENT Bloc (Ligne*,coordonnees)*>
<!-- a "Bloc" is made up of one or several "ligne" entries and their
corresponding coordinates -->
<ATTLIST Bloc Num ID #REQUIRED Type {corps_de_texte | marge |
has_de_page | has_de_page} #REQUIRED attribut {normal |
BIFPA | encadré} #REQUIRED
<!-- a "Bloc" entry has for attributes one ID "Num", a required "Type"
defined by: {body of text | margin | bottom of the page |
top of the page}, and a required attribut among {normal | BIFPA | squared -->
<ELEMENT Ligne (texte+)*>
<!-- a "Ligne" entry contains at least one "texte" element -->
<ATTLIST Ligne Type {Ligne | Inter-Ligne} #REQUIRED
<!-- a "Ligne" element has one mandatory attribut "Type" among {Ligne | Inter-Ligne} -->
<ELEMENT coordonnees (point,point)*>
<!-- coordinates are defined by a set of points -->
<ELEMENT texte (#PCDATA)*>
<!-- the textual ASCII entry of a text Bloc -->
<ATTLIST texte Num CHAYA #REQUIRED Type {Normal | barré
| souligné | illisible} #REQUIRED
<!-- a "texte" element has for attributes one mandatory id "Num" and one
mandatory attribut "Type" among {Normal | barré | souligné | illisible} -->
<ELEMENT point (#PCDATA)*>
<!-- the "point" element contains the point coordinates (x,y) -->
```

IA

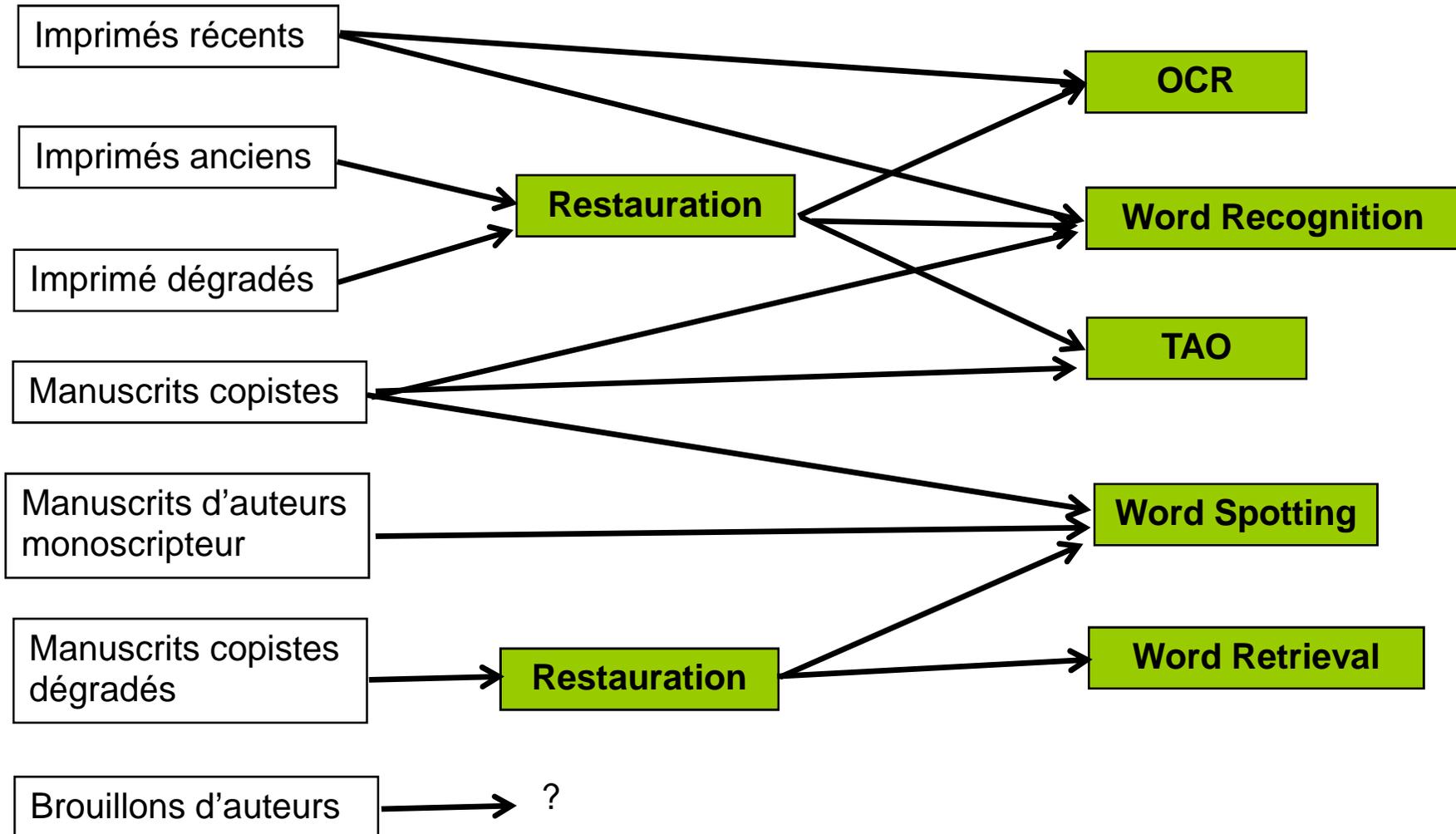
Post-traitements
Cohérence logique,
Linguistique...



Brouillons d'auteurs → ?

Comment accéder aux textes des documents numérisés ?

Cela dépend des images et de leurs contenus

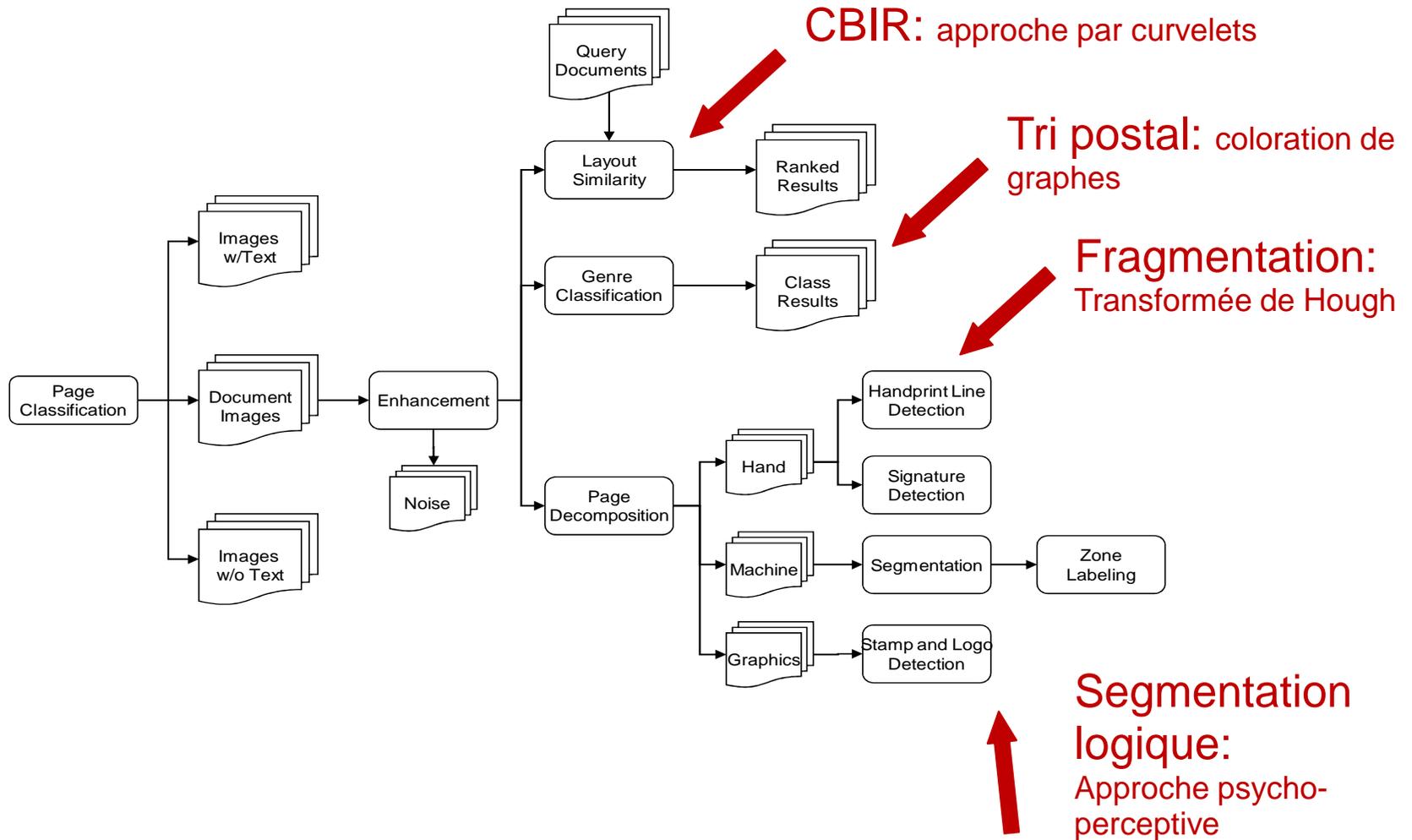


Partie 1

L'ANALYSE DE DOCUMENT



Analyse de documents



Une grande variété de tâches pour une grande variété de contenus

Quelles sont les difficultés inhérentes à l'analyse de documents ?

- .Tout est représenté dans un tableau "2D"
- .Tout est symbole
- .Les variations de symboles sont très nombreuses: épaisseur, taille, orientation...



.... et lorsque l'on a affaire à du manuscrit?

Sur les documents réguliers composites, on sait assez bien faire

- .Correction de l'inclinaison

 - *Orientation primaire des lignes de texte*

- . Détection des images et des régions de texte

 - Texture et orientation dominantes

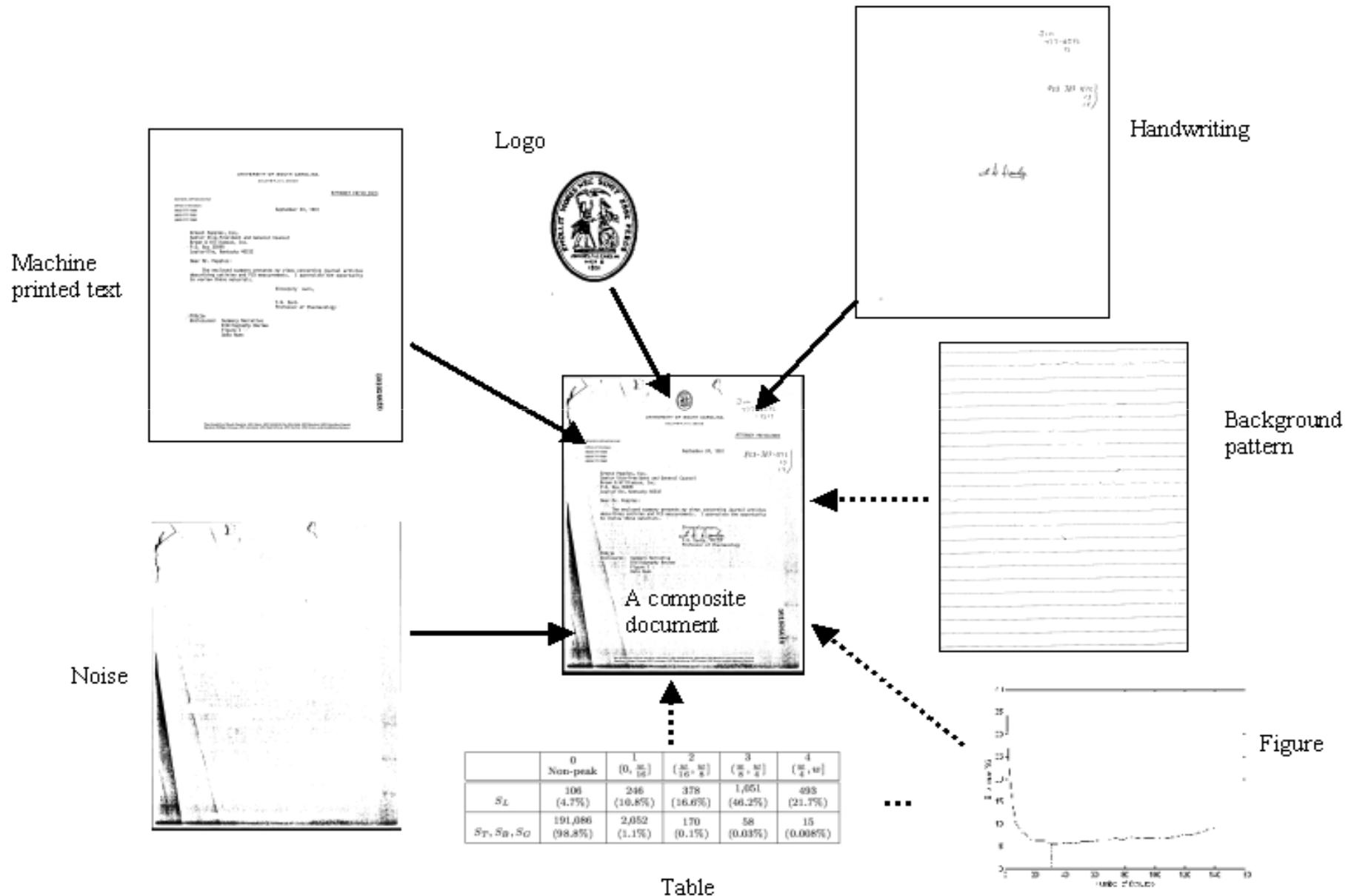
- .Classification structurelle

 - Inférence de la structure physique sur la logique

- .Classification des région de texte

 - Titre, auteur, notes entêtes, signature, etc.

Sur les documents composites bruités?



Principe général de la segmentation des documents de structures régulières

.Typiquement basé sur l'analyse des proximités spatiales

- Plages blanches
- Marges
- Différences dans les contenus

.Peut être très sensible au bruit

.Distinction entre deux types d'approches

- Top Down – du global au local
- Bottom up – du local au global

Sur les documents de structure irrégulière de contenu hétérogène?

- .On peut tenter des approches sur des régions localisées
- .Il reste à localiser les régions !!
- .On peut essayer de classifier les contenus: faut il parvenir à les caractériser initialement ?
- .La tâche est rendue difficile, parfois impossible...

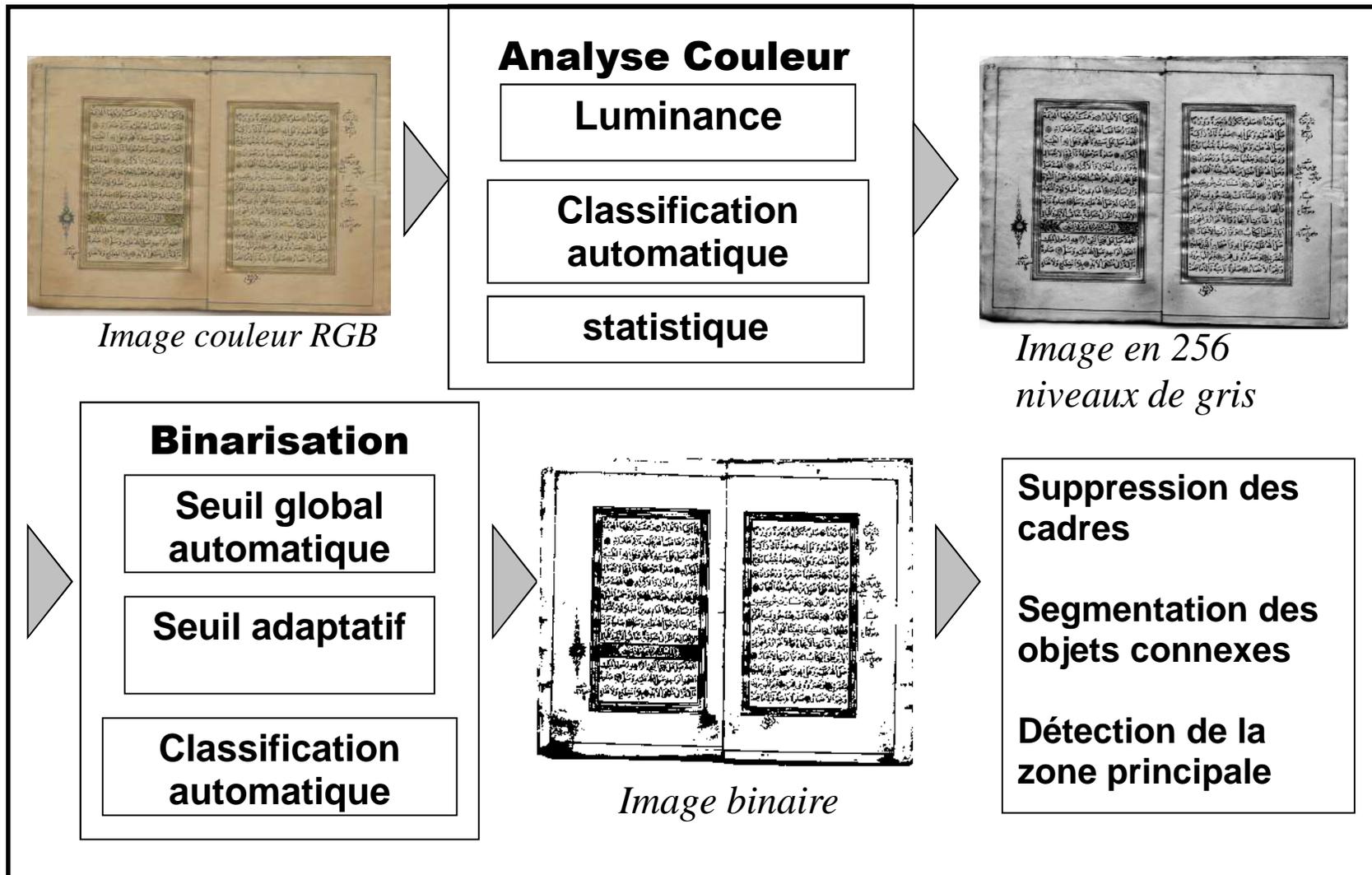


De la structure physique à la structure logique (ordre de lecture)

- .Difficile à obtenir selon le cas: pas toujours stable
- .Existence de guides explicites
 - “Continue à la page 4”
- .Existences d’indices structurels
 - Continuité d’un texte sur deux colonnes
- .Analyse du contenu sémantique
 - Statistique d’occurrences de mots, analyse de syntaxe

SEGMENTATION PAR L'EXEMPLE

Architecture d'un système : Segmentation



Segmentation

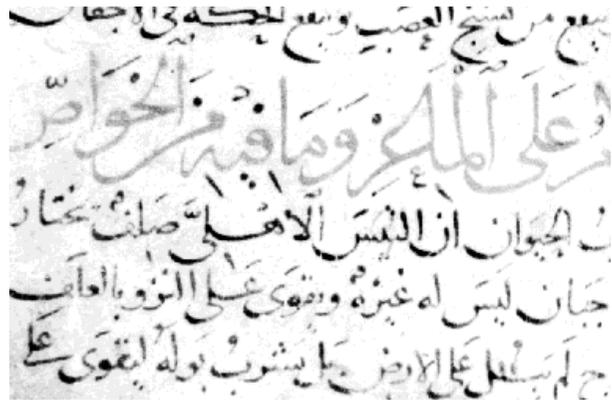
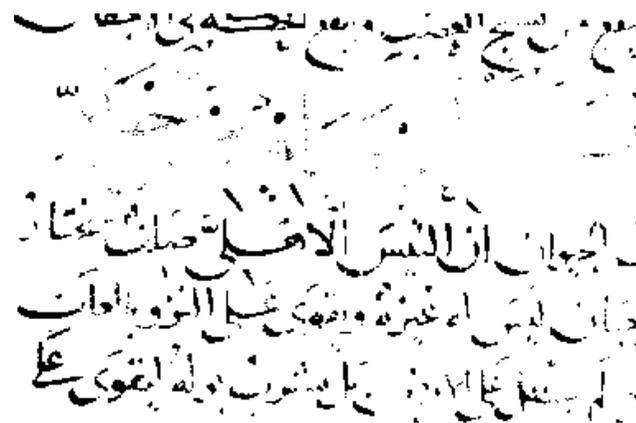


Image originale 16 niveaux de gris



Seuillage automatique globale



Seuillage automatique adaptatif



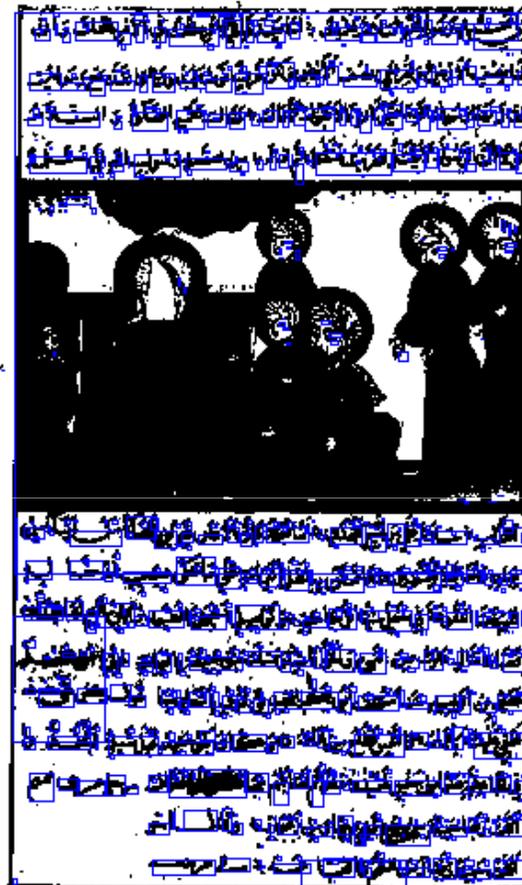
Seuillage par classification automatique

Recherche de tous les objets

الْتَهَامِي الْمَاضِي عَدْلِكَ لَكَانَ حُكْمُهُ دَعْوَى سَادِحَةٍ يُعْلَلُ بِهَا
طَرَفِي التَّقْيِيضِ الْمَخْتَصِرِ إِخْدِيَةً وَأَوْ يَتَنَبَّأُ بِالْآخِرِ وَقَدْ أَتَى عِنْدَنَا الْفَا
لِي الْفِعْلِ حَقٌّ يُوجَدُ وَالْمَاضِي مِنَ الْحُرُكَاتِ وَالْأَدْوَارِ وَالْآنَ
فِي الْعِلْمِ فَلَيْسَتْ بِالْإِنْفَاءِ وَهَذِهِ النُّقْطَةُ بِمَا يَكْفِي بِهَا

Segmentation des objets connexes

Problèmes des cadres et lignes

A page from an Arabic manuscript. At the top, there is a circular emblem or seal. Below it, there are several lines of text in Arabic script, with small blue boxes highlighting specific words or phrases. The bottom half of the page is dominated by a large grid table. The table has multiple columns and rows, with text in Arabic script filling the cells. The grid is outlined in black, and the text is in black.

Les cadres et tableaux forment aussi des objets connexes

Suppression automatique des cadres



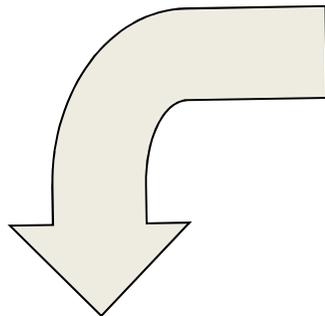
Pour simplifier l'image, il faut retirer les cadres, bordures et bords des ouvrages

Mesure le diamètre vertical maximal des objets

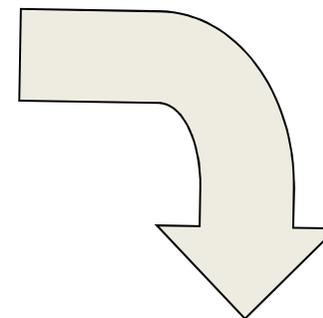


Chaque objet est rempli d'un niveau de gris correspondant à son diamètre maximal en hauteur

Classification Texte/graphique



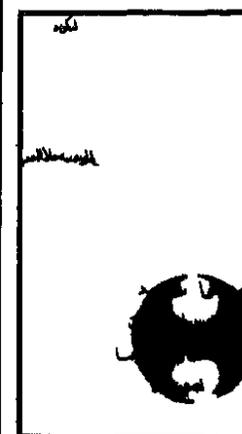
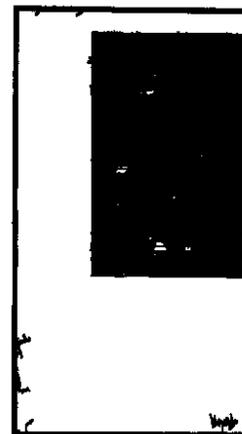
Objets
de faible diamètre



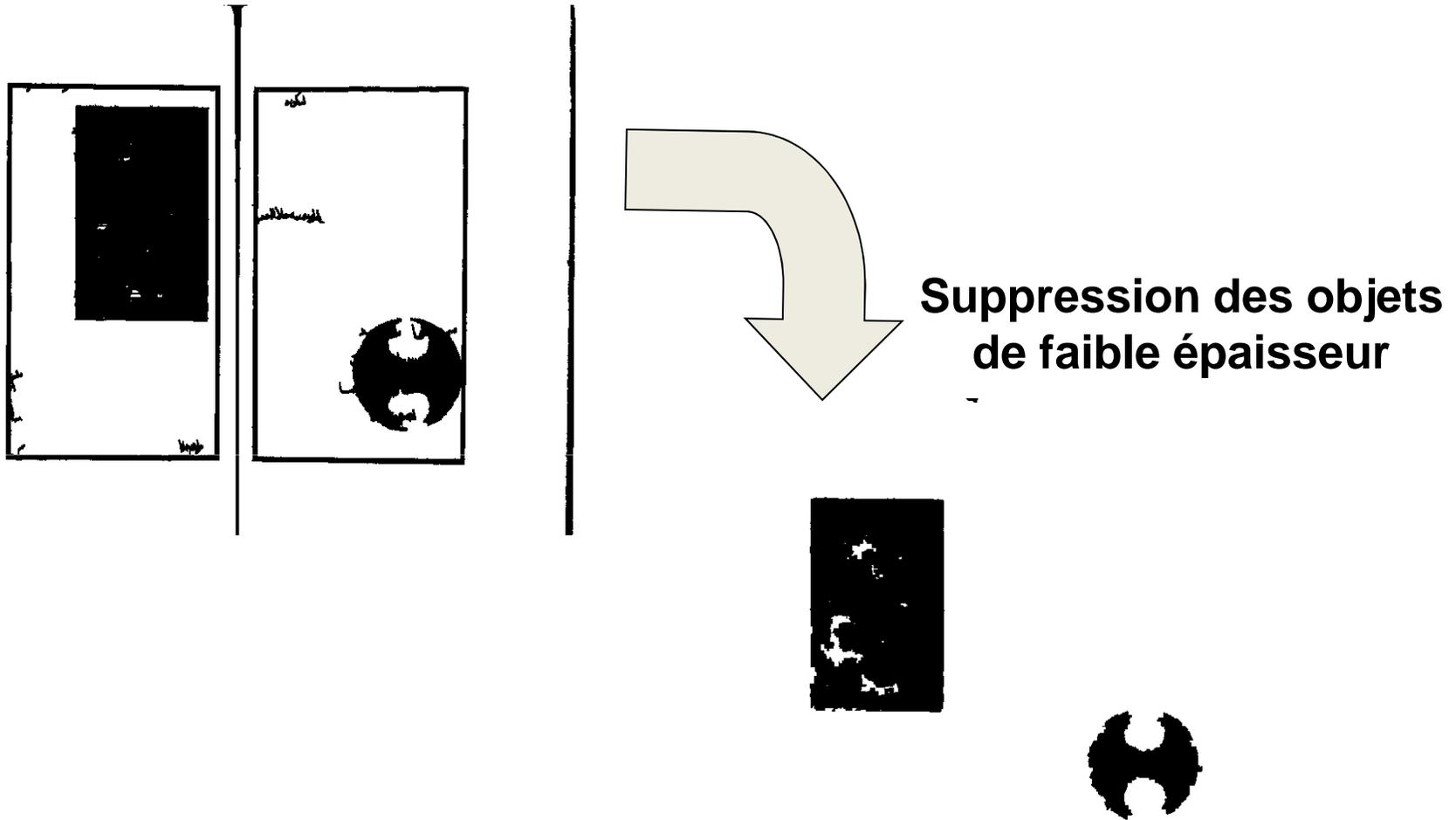
Objets de
diamètre élevé

أذا كنت في
التي من الرابع
ذا الركة، وتلك
بواسطة أو من
بواسطة لا تتدوم
وأولها إلى الأخر
القرآن تسمى
المؤلفين وتنفذها
بوزن التيمم الثلاثة أو الشاسع والعشرون من وزن الجوز
جميع الزوائد الخمسة من المعدن المنكسر من سلكها
حذره من سبب وأدوية من الشوك والملاط والملاط
المعبر شيا عليها إلى أن يكون أحدهم في بعض الزوائد
بوزنها في المؤمل تحت الشوك المائل إلى أن يكون

الذهب كالأرغمة وتنفذها في الأوزن الثلاثة أو الشاسع
من وزن الأربعة من المعدن المنكسر من سلكها
بوزنها من المعدن المنكسر من سلكها



Suppression des objets de faible épaisseur



Partie 3

OCR

ET ALTERNATIVES PAR TAO ET WS

Optical Character Recognition

- Documents imprimés
- Segmentation
 - Blocs
 - Colonnes
 - Paragraphes...
 - Mots
 - Caractères
- Reconnaissance
 - Caractères
 - Guidée par le contexte
 - Probabilité d'apparition
 - Dictionnaire
- Vérification
 - Dictionnaire

Optical Character Recognition

.Les approches usuelles

.Approche type « Pattern-matching »

- *Standard approach in commercial systems*
- *Segment individual characters*
- *Recognize using a neural network classifier*

.Modèle des chaînes de Markov cachées

- *Segment into sub-character slices*
- *Limited lookahead to find best character choice*
- *Useful for connected scripts (e.g., Arabic)*

OCR Accuracy Problems

- .Character segmentation errors

 - In latin scripts, segmentation often changes “m” to “rn”

- .Character confusion

 - Characters with similar shapes often confounded

- .OCR on copies is much worse than on originals

 - Pixel bloom, character splitting, binding bend

- .Uncommon fonts can cause problems

 - If not used to train a neural network

Améliorer la précision des OCR

- Prétraiter les images

- Mathematical morphology for bloom and splitting
- Particularly important for degraded images

- Mettre les OCR en concurrences, les sérialiser

- Individual systems depend on specific training data

- Employer des correcteurs linguistiques

- Use confusion statistics, word lists, syntax, ...
- But more harmful errors might be introduced

Les limites des OCRs sur le 16ème

bat ille ihesus: q̄ quom̄ p̄mū aufes uocarēt moises figurā
ihesum uocari: ut dux militiæ delectus esset aduersus am
nabant filios israhel: et aduersariū debellaret p̄ noīs figu

esse sensum semitas queritur. tanq̄ illi ad cogita
quadrigis opus eēt. Democritus quasi in puteo q̄
ut fundus sit nullus: ueritatem iacere demersam

RVDIMENTA. 6
hoc sanctum, sanctius, sanctissimum. Vocatiuo ô sãcte
sanctior, sanctissime, ô sancta, sanctior, sanctissima, ô
sanctum, sanctius, sanctissimum. Ablatiuo ab hoc san-
cto, sanctiõre vel sanctiõri, sanctissimo, ab hac sancta,
sanctiõre vel sanctiõri, sanctissima, ab hoc sancto, san-
ctiõre vel sanctiõri, sanctissimo. Et pluraliter. Nomina-
tiuio hi sancti, sãctiores, sanctissimi, hæ sanctæ, sanctiõ-
res, sanctissimæ, hæc sancta, sanctiora, sanctissima. Geni-

p̄ellent diei ac nocti: ⁊ diuideret lucem
ac tenebras. Et uidit d̄ q̄ esset bonū;
et factū ē uespere et mane dies quart⁹.
Dixit etiam deus. Producant aque
reptile anime uiuentis et uolacile sup
terram: sub firmamēto celi. Creauitq;
deus cete grandia. et omnē animā ui-
uentem atq; motabilem quā produxe-
rant aque in species suas: ⁊ omne vo-

ὅς ἀνοίη σέ, δικαιοῦς ἴδι σιωπῆ.
ἑτέροις πίδου πλέοι ἢ σιωπῆ. ὡς
Quand on te loue en vertu ou science,
Iuge cela toy mesme en conscience.
im aliquis, Quand aucun
lat te, te loue, te donne louange, dit bien
nénto. pense ⁊ sois aduerti

force des verbes substantifs ⁊
verbes passifs ⁊ les neutres, ⁊
mens, ⁊ les communs: comme,
o rectus. Scribo sedens. Per-
iniustus.
e vocatiui data neutris,
gnatio sæpe.

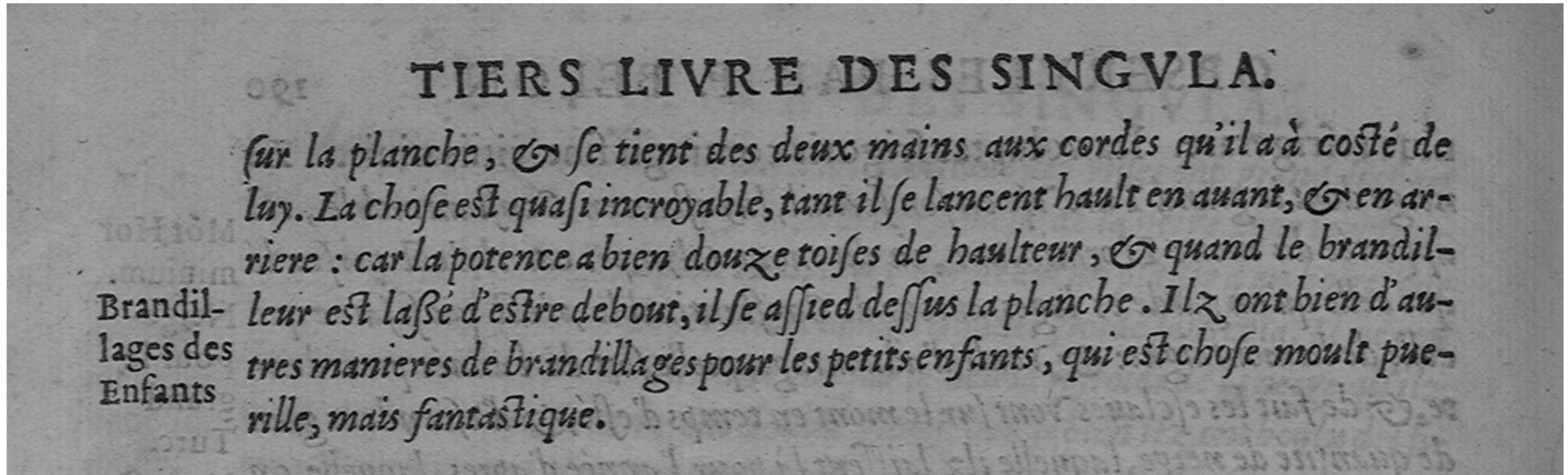
Polices et styles obsolètes que les OCRs ne reconnaissent pas

Et Œ st sp st fi æ ij ff ꝑ e

Caractères ligaturés aujourd'hui obsolètes

problème de codage UNICODE des caractères de la Renaissance

Imprimé de la Renaissance



Sortie OCR

TIERS HVRE „Es Si NG VL A~

		tient des. deùx; thé	ST 4UX: Cø? dt5	F COŠ16 de
		lu, ... Mb. ofieffquafiiiwiô414	taftt-Itft læfiççfli hatilt en	ariant, 4?
		rière~ :t4tL4pOtCfttC4biç%âØU~~tO~fts	dehaulteur.	quans w brandli-
ira~hdS1a		leur th læßé d'ëftn debost, iljc affiedde]	u~ laPI4ñChC	•Äilt~ ont ken .d'ij~.
lagçš çks		tres mænie Tes: dc brandillages pour les peurs	Èr?flmnu; 1m 4.	st chose~moufr puca~
Enfants		ñl14 mats		

A C E.

A C È, f. f. *Ace*. Ville de Phénicie, dans Strabon & dans Étienne. Ce fut depuis Ptolémaïs. Voyez ce mot.

A C É E, f. f. Ce mot se disoit autrefois pour *bécasse* : il vient d'*aceia*, qui vient d'*acus* à cause du long bec de la *bécasse*.

A C E P H A L E, f. m. *Acephalus*. Proprement, qui n'a point de Chef, de l'*a* privatif, & de κεφαλή tête, Chef. On a donné ce nom, 1^o à ceux qui dans l'affaire du Concile d'Éphèse ne voulurent suivre ni S. Cyrille, ni S. Jean d'Antioche : 2^o à des Hérétiques du V. siècle qui suivirent d'abord Pierre Mongus, ou Moggus : puis l'abandonnèrent, parce qu'il souscrivit au Concile de Chalcédoine. Ils suivoient les erreurs d'Eutichès. Et sous l'Empire de Justin, les Sectateurs de Sévère d'Antioche, & généralement tous ceux qui ne voulurent pas recevoir le Concile de Chalcédoine, furent appellez *Acephales*. Quelques uns prétendent que ce nom signifie *hésitant*, & que parce qu'ils tenoient la neutralité pour les décrets du Concile de Chalcédoine,

Résultat de TexteBridge (Xerox)

A C L

- A C ~, C f. *Ace*. Vill~ d~ Øhénicie~ dans Stra~, ou & dans t-tiei~ne. Ce fut depuis Ptôl~maVs. Voyez ce mot.
- A C _~. E, fi f. Cé mot se dif~ir autrefois pour *bJCAJJ'e* : il vient d'ac~ qui vient d'*acus* ~i c~ufe du long Lec *delab~4ffe* C E P H A L E, fi n. *Acephalus*& Proprement, qui n'a point de Chef, de l'a privatif, & de ~u4tête, Chéf. On ~i donné ce nort, 10~ ceux qui dans t'affaire du Concile 4~Éphèse ne votl— lurent fuivre ni S. Cyrille, ni S. Jean d'Antioche 2 z^o ~ desH~r&. tiques ~lu V~ fi~cle qui fuivirent d'abord Pierre Mongus, oU Moggiis~ ~ùis i'abandonu~renr, parce qu'il foufcrivit au Conde de Çhalc~doi~ie, Ils fuivQient les erreurs d'~utich~s. Et fous l'Empire de Ju<tin, les S~ateurs de S~vét~ d'Ant~oche, & g6-4~éralement tous ceux qui ne voplurent pa~ recevo~de Concile de Ch~lédoin&, furent appdlle~ *Acephales*.; Quelque~ uns pr&~ndent ~uec~noni fwnifie *héitant*, & que parce qu'il~te~ hoient la neutralité pour l'es decr~rs dit Concile de Ch~lcédoine~ qu'ils ne *fed&èrru~nd~entd* rien qu'ils h~fitoient quancipn les préflbit, ~1S furent appellez *Acephales* c³ea~dire, *hóftatis*. Mais Vautre opinion ea ph~ vraie, & *Ac~pbale* n'~ point ce fens. Voyez Bolland. T~ I. Anaaafùs le ~ibltothéca~r~ appelle cette ~emption de la frdfd~aion ~u P. itiarche, Atinocéphalie, ~*autocephal~a*; ; on a ~appel l~ *Acephat~s* le~ Clercs qw ne vivoient pàs fous'la Difcipline Éccí~fiatUqÙe d'un ~v~que Afidore *PcEcclef off. La; 11J*. Les Conciles de Majtence, Can. 22. de Meaux I an~~y. Can.~. de PaAsCan~i o~ de Pavie cri 8yo. Can~i ~. &c. ont fait différens r~glernens c~n~re ces *ClercsAcephales*. 'On *cii* trônve ericore dans tes Capitul aires de Ch~trics le Chais. Ve, **L. VI. C. 57**. dari~ lurchard, **L. II. C. 26**. dans kegin~n ~ l'an de Jefus-Chrift S. ~ Baronius ~ U~ntiée i o~o~ ~iucbert frère de Thi~rb~rge Concubine de. Î~oçhaire, fut appellé *4ce~ eha7é*, parce q~xe <coinrne di~n(1~ Araiales de Mets ûl~

Résultat de Omnipage (Caere)

A C 1:

AC , Γ f. Àce. *Ville de phénicie, dans ~trabdn & dans 4'tierine.*

Ce fut depuis Ptôlémaïs. Voyez ce mot:

r C \$ E, f f. Cè tnot fe droit autrefois pour *hécasse* : *il vient d'ac= cela, qui vient d'acus à esuse du long bec delahécaTe.*

â~ C E P H A L E, fin. *Acephalusi* Proprement, qui n'a point de Chef, de l'a privatif, & de μαρ, ~À~ tête, Chèç, 011,1 donné ce 11011, i ° à ceux qui dans l'affaire du Concile 'Éphète ne voulurent fuivre ni S. Cyrille, ni S. Jean d'Antioche' ~ z' à des Hçrés~: tiques du V; siécle qui fuivirent d'abord Pierre Mongus, ou 1Vloggus: puis sabandonnèrent, parce qu'il toufcrit au Concile de Glzalcédoiaè. Ils fuivoient les erreurs d'putichés. \$t sous l'Empire de justin, les Sédateurs de Sévère d'Antioche , & gé'héalement tous ceux qui ne voulurent pas recevoir le Concile de Châltédoinè, furent appellez *Acephatès*.: *Quelque, à uns prétendent que ce nom signifie h~ritant , & que parce qu'ils, tei*

Résultat de Finereader (Abby)

A C E .

- À C È, f. f. *Ace*. Ville de Phénicie, dans Strabon & dans Étieline. Ce fut depuis Ptolémaïs. ^ Voyez ce mot.
- A C E E, f. f. Ce mot se diloit autrefois pour *bécœffe* : il vient d'*acceia, qui vknt d'acusàcàute* du long bec à *ehèeca/e*.
- A C E P H A L E, f; m. *Acephalusi* Proprement, qui n'a point de Chef, de 1*4 privatif, & de *te^Ai* tête, Chef. On a donné ce nom, i° à ceux qui dans l'affaire du Concile d'Ephèse ne voulurent suivre ni S. Cyrille, ni S. Jean d'Antioché i 2, ° a' des Hérétiques du Vi siècle qui suivirent d'abord Pierre Mongus, ou Moggus S puis l'abandonnèrent, parce qu'il souscrivit au Concile de Chalcedoine, Ils suivirent les erreurs d'Éutychès. Et sous l'Empire de Justin, les Sectateurs de Sévère d'Antioche » & généralement tous ceux qui ne voulurent pas recevoir le Concile de Chalcedoine, furent appeliez *^cefiats*. Quelque^ uns prétendent que ce nom signifie *béfitant*, & que parce qu'ils se^ hoient la neutralité pour les décrets du Concile de Chalcedoine >

Transcription Assistée par Ordinateur (TAO)

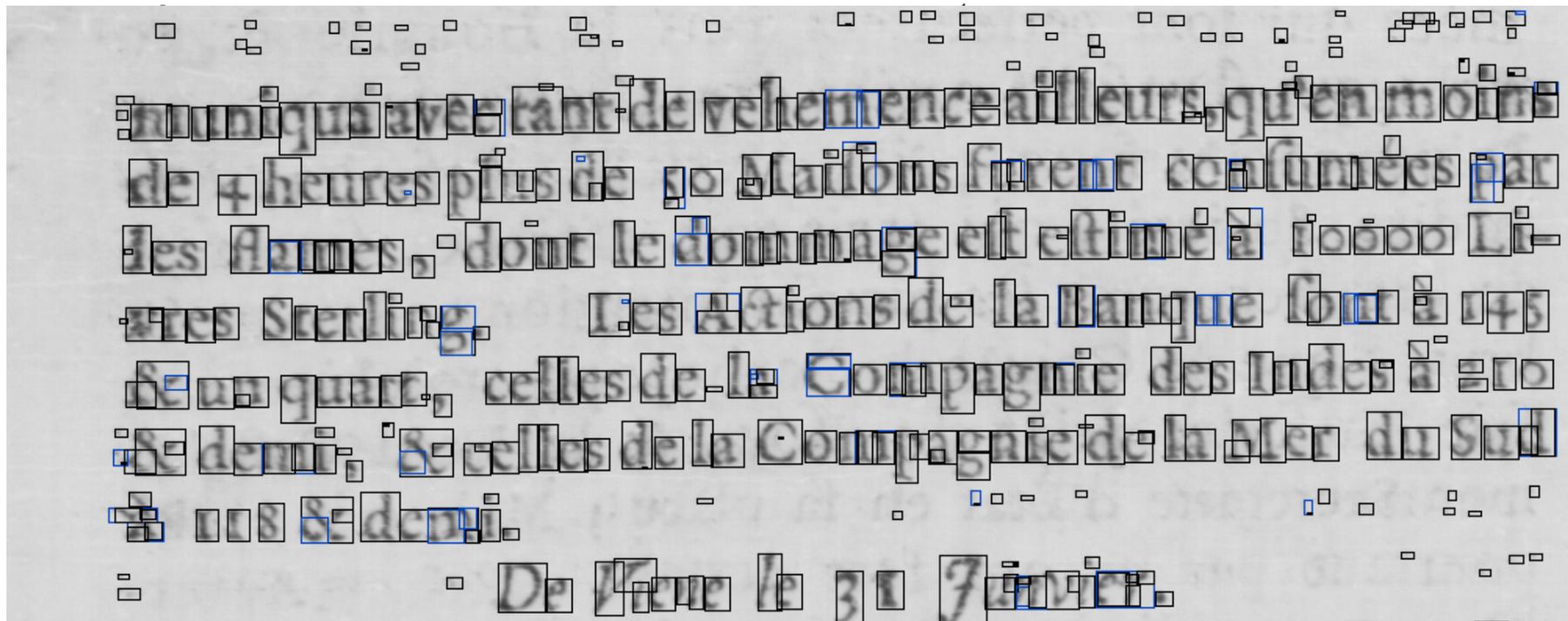
- *Documents imprimés*
- Segmentation
 - Lignes
 - Mots
 - Caractères
- Prototypage / Groupement (clustering)
 - Caractères
- Étiquetage manuel des groupes

Principes du prototypage

Nouvelle segmentation en caractères guidée par la redondance des formes basé sur les hypothèses suivantes :

H1 : pour un caractère fragmenté, il existe au moins une forme correctement imprimée du même caractère

H2 : pour deux caractères collés, il existe au moins une occurrence de ces caractères qui sont imprimés isolément.



Transcription Assistée par Ordinateur (TAO)

500 caractères

80 prototypes

For the Olympics, hundreds of Korean-Americans have flown to Seoul, according to managers of Korean Air. Joining them were at least 50 students, who are serving as translators at the games.

Other young Korean-Americans, who often have to translate English for their parents, enjoy the Olympic view from home. "It's great to see on television how beautiful it is and what it's like in Korea," said Julie Hong, 21 years old, who left Korea with her family when she was 4.

The Olympics were celebrated here two weeks ago with several floats in the annual Korea Days parade. A

123 456 78910 11121314, 51516173181714 219 202
21182223- 24 10 183125221614 5222518 262722816 42
2213132317313233 42 343532363318314 3719 2038361
24313. 3921240314041 4251810 2818318 364 31836
171832414, 2852 46 318 146325311633 2214 432216
434 4518 334610 1814.

74518 3 492304041 202364616- 24 10 631213473214,
2852 21941832 5475018 42 43363214846451 5216531
26221 54518123 1146 3340414, 183255249 4518 78
5631628 263210 5210 18. ' 574 14 33318364 582
58331825121412259 5228 501836305812193027 3142 6
124 14 8126118 1232 52236336, ' 14361264 653027
491847314 26954, 2852 2718194 20231846 281245 1
193610 12689 285640 1456 704714 71.

72736 727910 11121314 2818318 1318686503464674 |
4702 7518186114 46332 7576473 14650634627 19824
58518 46323230363 5223646 7746914 1146 3471718.

Transcription Assistée par Ordinateur (TAO)

The image displays the 'Transcripteur' software interface. On the left, a manuscript page is shown with a decorative header and a large initial 'M'. The text reads: 'A TRESILLVSTRE ET REVERENDIS-SIME SEIGNEVR, FRANCOIS CARDINAL DE TOVRNON, SINGulier & liberal Mecenas des hommes studieux de vertu, Pierre Belon son resumble done si que seruiteur salu, & entiere prosperite.' Below this, a paragraph begins with 'Onseigneur, c'est à bon droict que les gents doctes vous ont en admiration, & que le peup e-stranger affecte à nostre republi que, comme ausi le François a grandement loué & estimé l'ex-celleuce de vostre bon jugemêt, & magnifié vostre prudence & vertu: car entre tous autres illu-stres prelatz, vous avez singulierement aimé & honoré les lettres, aduacé les lettrez, & par vostre speciale faueur...'. The software window 'Transcripteur' is overlaid on the manuscript. It shows the text 'stus de latin bleu' with a search bar containing 'a'. Below the search bar, there are navigation buttons and a progress bar for '% de prototypes à saisir' and '% de texte transcrit'. On the right, a grid of characters is displayed, including letters, ligatures, and symbols, each with a small number indicating its position in the character set.

La saisie de 2% des caractères, réalisée en 6 heures, suffit à reconstruire un livre du XVI transcrit pour une recherche par mot

Transcription parfaite en quelques heures

- TAO implantée dans DEBORA (projet Européen 1998-2003)



Comment étudier les erreurs de la TAO ?

En affichant toutes les formes par prototype

The image displays a grid of character prototypes for TAO error analysis. The grid is organized into several rows of characters, with specific errors highlighted by callouts:

- Erreur de Substitution:** A blue callout points to a 'G' character in the first row, which is highlighted with a blue box. This row contains a sequence of 'G' characters, with the 10th character being the one pointed to.
- Etiquette Multiple:** A green callout points to a 'c' character in the second row, which is highlighted with a green box. This row contains a sequence of 'c' characters, with the 25th character being the one pointed to.

The grid also shows other rows of characters, including 'e', 'a', 'l', 'd', and 'u', which are used as prototypes for various TAO errors.

Transcription parfaite en quelques heures

The screenshot displays the DEBORA v0.1b10 software interface. On the left, a tree view shows the document structure for 'Magnificence', including metadata and a list of pages from 1 to 19. The main window shows 'Page: 3 (21%)' of a manuscript. The page features a large, ornate initial 'L' and text in French. A white callout box with the word 'Transcription' points to the text. Below the manuscript image, there are tabs for 'Texte brut', 'Texte lemmatisé', and 'Texte moderne'. The 'Texte moderne' tab is active, showing the transcribed text. On the right side, a thumbnail view shows a grid of pages, with 'Page 3' and 'Page 4' highlighted.

DEBORA v0.1b10
Fichier Edition Rechercher Affichage Fenêtres

Magnificence
Structure de l'ouvrage:
- Titre: Magnificence
- Auteur: COLLECTIF
- Bibliotheque: Bibliothèque Nationale c
- Lieu de publication: Lyon
- Date de publication: 1549
- Editeur: Rouillé
- Langue: Français
- Cote: Réserve 355882
- Sujet: Magnificence de la superbe et t
Pièces préliminaires:
+ Page: 1
+ Page: 2
Pages de Texte:
+ Page: 3
- Image composite:
- Imagette:
- Imagette:
+ Page: 4
+ Page: 5
- Image composite:
+ Page: 6
+ Page: 7
+ Page: 8
+ Page: 9
+ Page: 10
+ Page: 11
+ Page: 12
+ Page: 13
+ Page: 14
+ Page: 15
+ Page: 16
+ Page: 17
+ Page: 18
+ Page: 19

Titre: Magnificence .Page: 3 (21%)

de Monseigneur
d'or cliquant

Transcription

LE TRES
itien Roy de France He
deuxiesme uolant a
heureux aduenement u
les Frontieres de son
aulme, comme Prince
dent, delibera de pass
Piedmont pour y ue
forteresses, & pour plu
autres grandz respec
de la s'en retourner par Lyon. Ce que Monseigneur le Re
uerendissime Cardinal de Ferrare, Archeuesque & Conte de
Lyon, & Monseigneur le Gouverneur firent diligemment
entendre à Monfieur le Lieutenant du Roy, & Messieurs
les Coseilliers & Etcheuins de la Ville pour se preparer à le
recepuoir à son retour. Parquoy Messieurs de la Ville, ne
uolant degenerer à leur antique generosité Romaine, co-

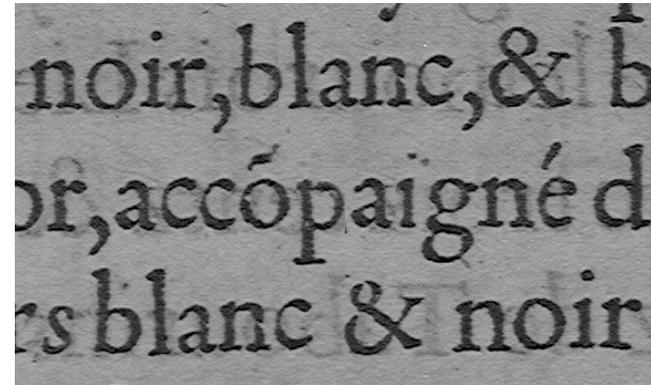
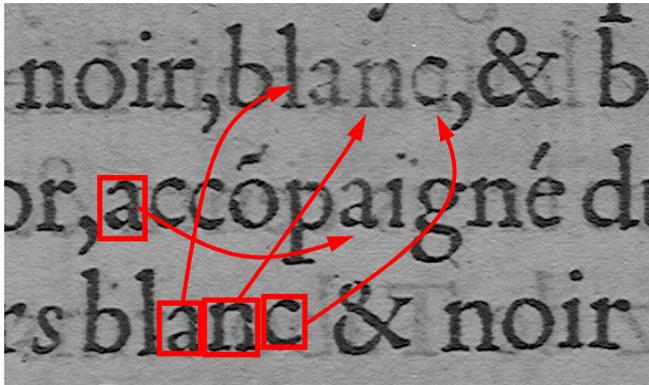
Texte brut | Texte lemmatisé | Texte moderne

de la sen retourner par Lyon Ce que Monseigneur le Re
uerendissime Cardinal de Ferrare, Archeuesque & Conte de
Lyon, & MMonseigneur le Gouverneur firent diligemment
entendre à Monfieur le Lieutenant du Roy, & Messieurs
les Cofeilliers & Etcheuins de la Ville pour se preparer à le
recepuoir a fon retour Parquoy Messieurs de la Ville, ne
uolant degenerer à leur antique generosité Romaine, co-

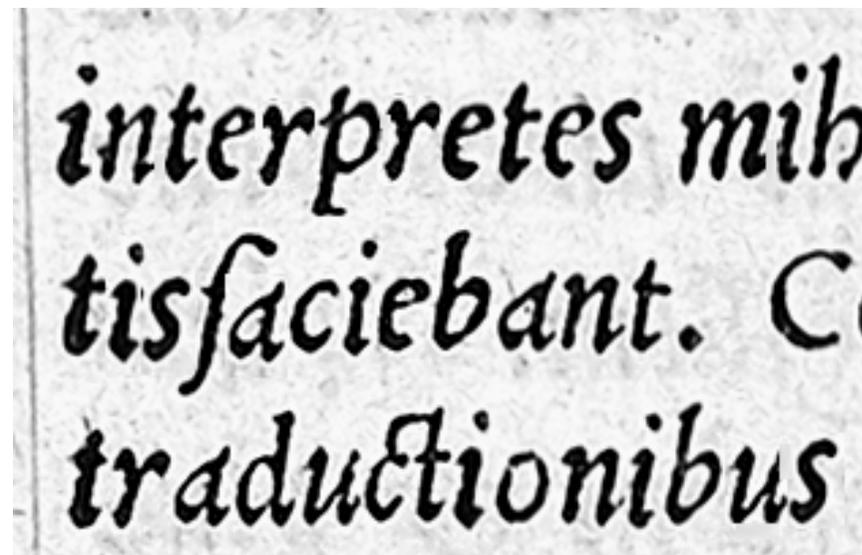
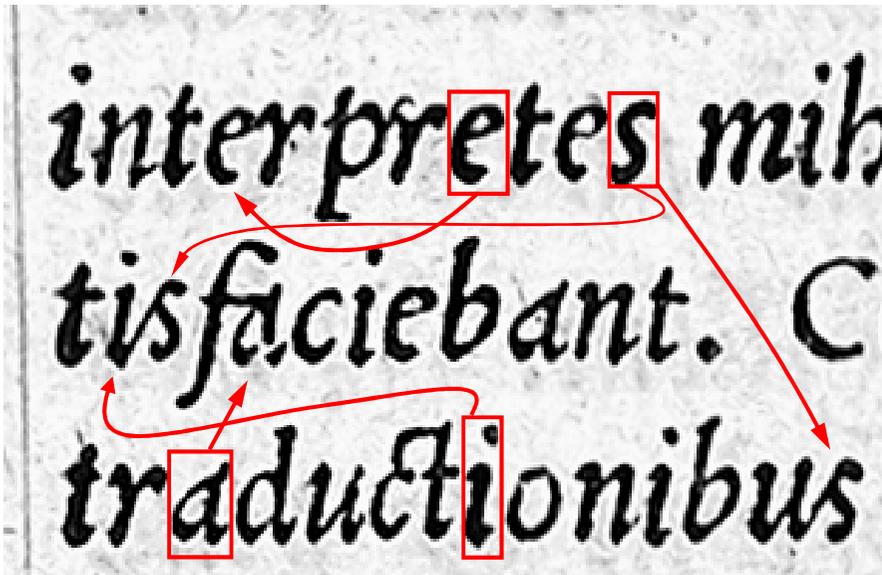
Page 3 Page 4
Page 8 Page 9
Page 13 Page 14

Principes de reconstruction

Phase 1 : Recouvrir les caractères cassés par la forme entière similaire.



Phase 2 : Recouvrir une forme collée par des formes isolées similaires



TAO: ouverture

- Compression

- Dictionnaire de formes
- Compression JPEG du fond (papier)

- Restauration

- Substitution des caractères par leur prototype

- Classement des mots par typographie

- Les polices et différentes typographies (gras, italiques) génèrent des prototypes différents

- Reconnaissance des mots par décodage

- Dictionnaire

Découvertes grâce à la TAO

1) Les typographes de la Renaissance composaient les mots fréquents avec les mêmes plombs

ant d'eulx leurs Ta
pres eulx cent foi.
dans eulx. Et cepei
apres eulx deux cét.

2) Une augmentation soudaine du nombre de nouveaux prototypes de caractères témoigne d'une « rupture » dans la fabrication comme l'insertion de nouvelles pages créées dans d'autres ateliers et nous renseigne sur la fabrication des livres

Résultat d'une recherche par mot dans les données et les métadonnées

The screenshot shows the DEBORA v0.1b7 application interface. The main window displays a document titled "Magnificence" with a search for the word "Lyon". A search results window is open, showing the following results:

- Texte> Bibliotheque Bibliothèque Nationale de Lyon
- Texte> Lieu de publication Lyon
- Texte> Sujet Magnificence de la superbe et triumpante entrée de la noble & antique cit
- Image> Magnificence Page n°2
- Image> Magnificence Page n°2
- Image> Magnificence Page n°2
- Image> Magnificence Page n°3
- Image> Magnificence Page n°3
- Image> Magnificence Page n°21

The main document window shows a page with a large illuminated initial 'P' and the word 'Privileg'. The search results are highlighted in pink in the document text, such as 'Lyon' and 'Lyon'.

Limites d'une recherche de mots dans une transcription non lemmatisé

Recherche de "harque" ramène « harquebouse » mais pas « hacquebute » ni « hacquebusiers ». Pas de recherche imprécises avec substitution de lettres ou de jokers !

The screenshot shows the DEBORA v0.1b10 software interface. The main window displays a page of text from a historical document. The text is as follows:

luy deux laquais ueffus de latin bleu. Apres luy les Hacque-
busiers de la ville de troys à troys en nombre de troys centz
trente huict habilliez de blanc & noir. A icauoir, le collet &
chauffes de melours noir chargez de boutons & fers d'or : le
pourpoint de latin blanc, & doubleure de chauffes de taffe-
tas blanc rayé d'or:chaſcun ſon mourrion doré avec le pen-
nache de blanc & noir femé de pailletes d'or:la harquebou-
ſe & le reſte des autres armes ſemblablement dorez: accom-
paignez de leur Enſeigne ayant au milieu les armes de la
Ville, une hacquebute au deffouz, avec leurs Tabourins &
Fiffres de meſme liuree, pour un ioyeux commencement de
leur ſuytte.

Au dos deſquelz ſuyuoit la ſecode Bande, au premier ranc
de laquelle (ſelon la delibération de l'ordre, que dict a eſté cy

The search results are shown in a list at the bottom of the window:

Texte brut | Texte lemmatisé | Texte moderne

claves de lours noir chargez de boutons et fers dor * led
pourpoint de latin blauc et doblere dechauffes de taffed
tas blanc ray dorclafcun foin innourrion doré avec le pen-d
mache de blinc & noir fene de pailletes dor:la harquebou-d
ſe et le reſte des autres aries ſeinblableinei t dorez: accomd
paignez de leur Eieigne ay aint au uieu les armes de lad
Ville, une hacqueiebute au difouz, aec leurs Tabourins etd

Annotations in the image:

- A box labeled "Mots non retrouvés par le logiciel de recherche" points to the word "Hacquebusiers" in the text.
- A box labeled "Mot retrouvé par le logiciel" points to the word "harquebouse" in the text, which is highlighted in pink in the original image.

Recherche de mot-image (WS)

Sur les documents non transcriposables

- Proposée par la communauté de reconnaissance vocale « Tohlicek et al, 1989 », « Paul et Rause, 1990 ».
- Utilisée comme une alternative à l'OCR
 - Pour des expériences sur l'indexation des documents manuscrits « Kuo et Agazzi, 1994 », « Manmatha et al, 1996 », « Rath et Manmatha,2003 » ,« Rath et Manmatha, 2007 ».
 - Pour l'extraction d'informations des documents imprimés « chen et al, 1993 », « Chen et al, 1995 », « Yue et al,2004 », « Jawahar et al, 2004 » et manuscrits « Zhang et al, 2004 », « Srihari et al, 2005 », « Srihari et al, 2006 ».
- Plus efficace que l'OCR sur les documents dégradés et manuscrits.

Recherche de mot-image (WS)

Sur les documents non transcriposables

- Principe général:
 - Prétraitement:
 - Enlever le bruit, corriger l'inclinaison.
 - Dépend de la qualité du document sur lequel on travaille
 - Analyse (segmentation et extraction de caractéristiques):
 - Etape de segmentation du document:
 - Segmentation du document en lignes.
 - Segmentation du document en mots.
 - Segmentation du documents en mots puis en caractères.
 - Extraction des caractéristiques:
 - Chaque mot est décrit par un vecteur de caractéristiques ou une description structurelle.
 - Appariements:
 - Avec/ sans apprentissage
 - Appariement d'images (similarité, métrique)

Recherche de mot-image (WS)

Sur les documents non transcriposables

- Les techniques de segmentation se divisent en deux groupes.



démonstration

SANS SEGMENTATION: Le document tout entier est considéré comme une seule image et aucune segmentation n'est appliquée (Leydier et al, 2005), (Leydier et al, 2007).

- Les images de documents sont analysées globalement par une recherche de zones d'intérêt, de guide, d'éléments structurants: SIFT, Gradient cumulés, orientations par Curvelets...) puis une recherche séquentielle est produite.

AVEC SEGMENTATION: le document est considéré comme une collection de sous-images (mots et caractères).

- Les documents sont segmentés en sous images (CCs).
- Méthodes top-down vs bottom up
- Méthodes globales (holistiques: ligne et mot) vs locales (analytiques: caractère)
- Appariement séquentiel entre mot –requête et mot(s)-image

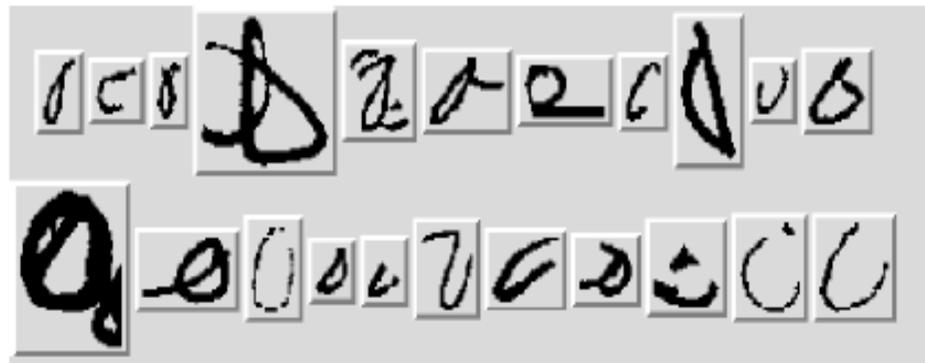
Recherche de mot-image (WS)

Sur les documents non transcriptibles

- Le mot vu dans sa globalité (*Madhvanath et al, 2001, Lavrenko et al, 2004, Madhvanath et Govindaraju, 2001*)



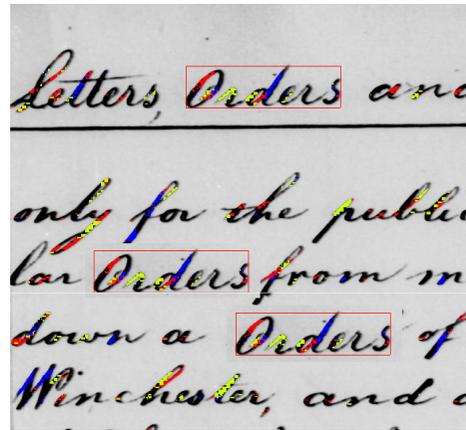
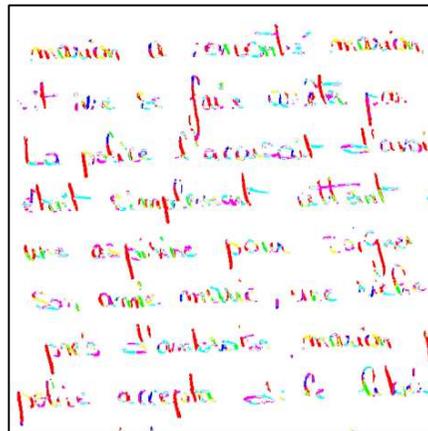
- Le mot vu comme une collection de sous unités (utilisation d'un lexique), *Gatos et al, 1997, Madhvanath et al, 2001*.



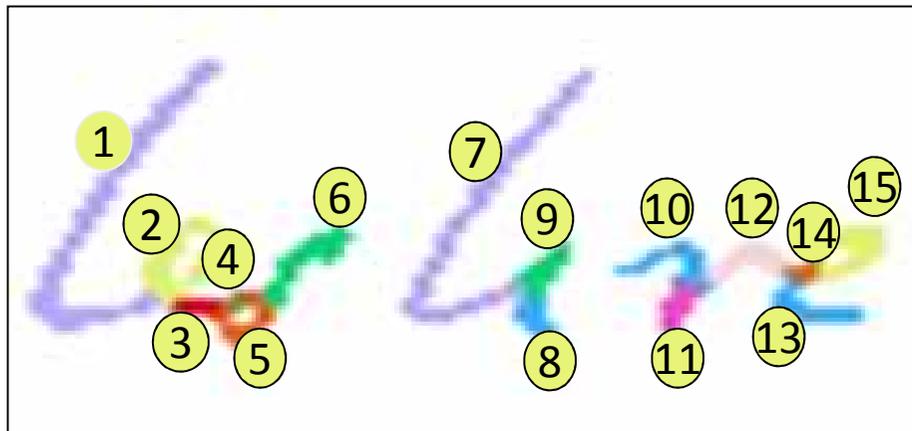
Recherche de mot-image (WS)

Sur les documents non transcriposables

- Le mot vu comme une collection d'orientations (*Joutel-eglin, 2008*)



- Le mot vu comme un ensemble de graphèmes, (Tayeb, Daher 2010).



Recherche de mot-image (WS)

Sur les documents non transcriposables

Techniques:

- Word Shape/HMM - (Chen et al, 1995)
- Word Image Matching - (Trenkle and Vogt, 1993; Hull et al)
- Character Stroke Features - (Decurtins and Chen, 1995)
- ▣Shape Coding - (Tanaka and Torii; Spitz 1995; Kia, 1996)
- ▣Gradient cumulés (Leydier, 2006)
- ▣Code Book (Tayeb 2007, Daher 2010)

Applications:

- ↻Navigation rapide de documents non transcrits (pas de version "texte- ASCII" disponible)
- ↻Accès aux documents manuscrits

Evaluation de systèmes :

- ↻Gribouillis et OCR (DeCurtins, SDIUT 1997, Yue et al, 2004): comparaison système de reconnaissance de mots avec système OCR, le système de word spotting est plus robuste.

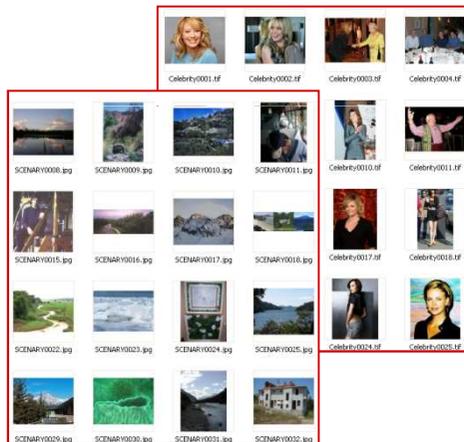
Partie 3

RECHERCHE D'INFORMATION

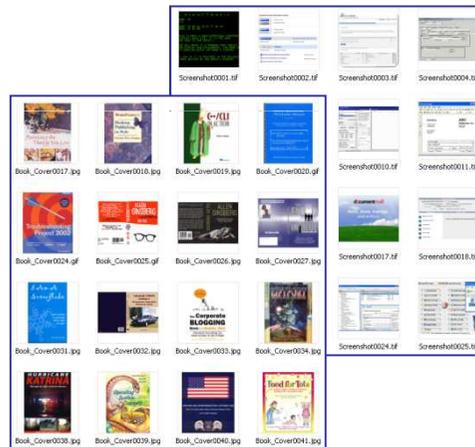
***CLASSIFICATION, INDEXATION, SIMILARITÉ
APPARIEMENT ET RECHERCHE D'INFORMATION***

Classification d'images par similarité

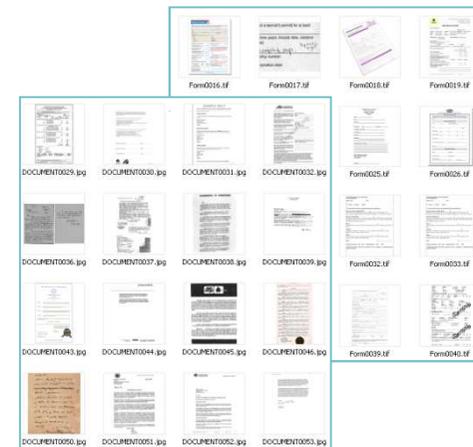
- Une grande variété de contenus
- Distinction: documents contenant ou pas du texte



Images



Images et texte



Images de Document

Indexer les images selon quels critères?

.La structure: indicateur important de contenu

-Courriers d'entreprises, documents administratifs, financiers

-Journaux, sommaires...

.Présence d'informations: connaissances a priori des contenus

-Présence de manuscrit?

-Bruit, résolution spécifique?

-Couleurs, niveaux de gris?

-Informations implicites de structures?

-Quelles caractéristiques discriminantes?

Quelles caractéristiques ?

Images de traits

- .Accumulation de motifs et détails structurés***
- .Géométrie, topologie, résolution, échelle...***
- .Des bandes de fréquences ciblées (en particulier, images de texte)***
- .Une information plus localisée liée au tracé***
- .Séparation forte forme / fond (information / environnement)***
- .Prédominance du segment (1D)***
- .Intention de l'auteur***

Images naturelles

- Couleurs ou niveaux de gris répartis avec cohérence spatiale sur des voisinages locaux***
- Couleur, texture, forme...***
- Accumulation de fréquences équiprobables***
- Information plus diffuse***
- Tout peut être potentiellement information***
- Prédominance de la surface (2D)***

Quelles caractérisations ?

.Dépendantes de l'objectif : *CBIR, catégorisation, restauration*

.Images de traits

- ≈ *Caractérisation locale de formes élémentaires (lettres, graphèmes, segments, divers traits...)*
- ≈ *Agencement global des formes locales (textures par textons)*
- ≈ *Caractérisation par blocs*
- ≈ *Forte composante géométrique*
- ≈ *Importance des détails, de la résolution ...*
- ≈ *Invariance à la rotation selon contexte*

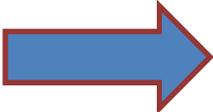
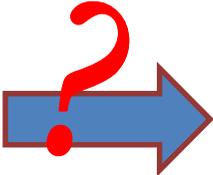
.Images naturelles

- ≈ *Informations locales, globales, par régions ...*
- ≈ *Caractérisation fréquentielle, ondelettes*
- ≈ *Information plus diffuse , souvent moins sensible à la résolution*
- ≈ *Invariance à la rotation très fréquente*

Approches et objectifs

	Locales	Par régions ou spécifiques	Globales
	Analyse de l'image au niveau pixel ou voisinage	Agglomération de caractéristiques locales Analyse par région	Analyse statistique de la multitude Transformées de l'image entière
Images de Traits	Analyse des contrastes locaux, de points caractéristiques...	Approches par régions d'intérêt, et fenêtre glissante Analyse de la géométrie, squelette WordSpotting, logos, LBA...	Agglomération de critères bas niveau, CodeBook de motifs, Fréq., Ondelettes, Curvelets, Similarité de scripteurs, catégorisation de documents (mise en page)
Images Naturelles	Points d'intérêt, descripteurs et invariants locaux...	Recherche d'objets, Contours (CSS), pattern matching... Visages, bâtiments, plaques...	Analyse fréquentielle, ondelettes, statistiques, texture... Classification sémantique, catégorisation d'images...

Quelques pistes ...

- Images naturelles  Images de traits
 - .Extraction de points d'intérêts type SIFT (Word Spotting)
 - .Caractéristiques Curvelets, Hermite (Reco. Scripteurs / Classe)
 - .Analyse de texture / Texte (Reco. Scripteurs / Classe)
- Images de traits  Images naturelles
 - . Analyse par composantes connexes élémentaires (Word Spotting, Reco. Scripteurs / Classe)

Mais le passage de l'un à l'autre n'est pas systématique !!! Il faut s'adapter au contexte qui constitue un environnement contraignant et impose de nombreuses limitations

QUELQUES EXEMPLES

Appariement de formes



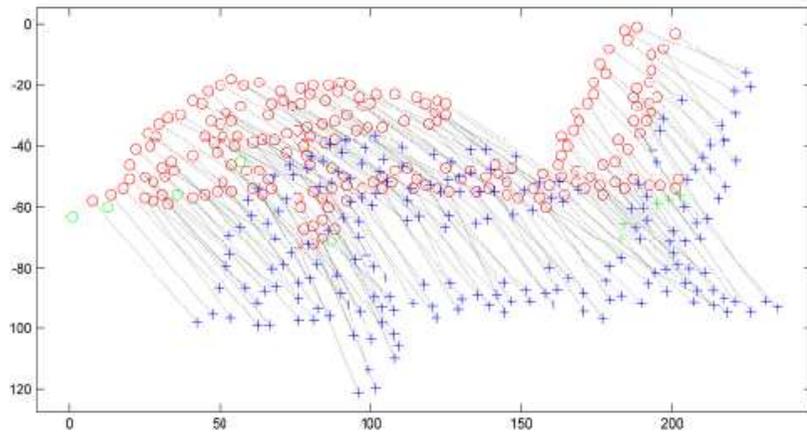
(a)

(b)

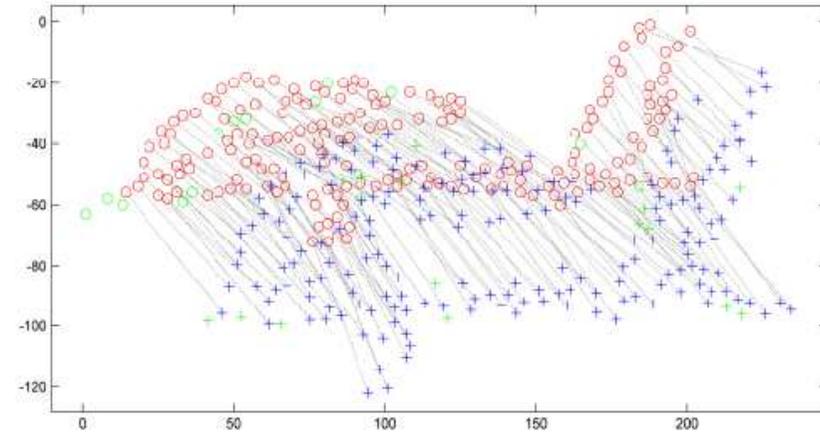


(d)

(e)



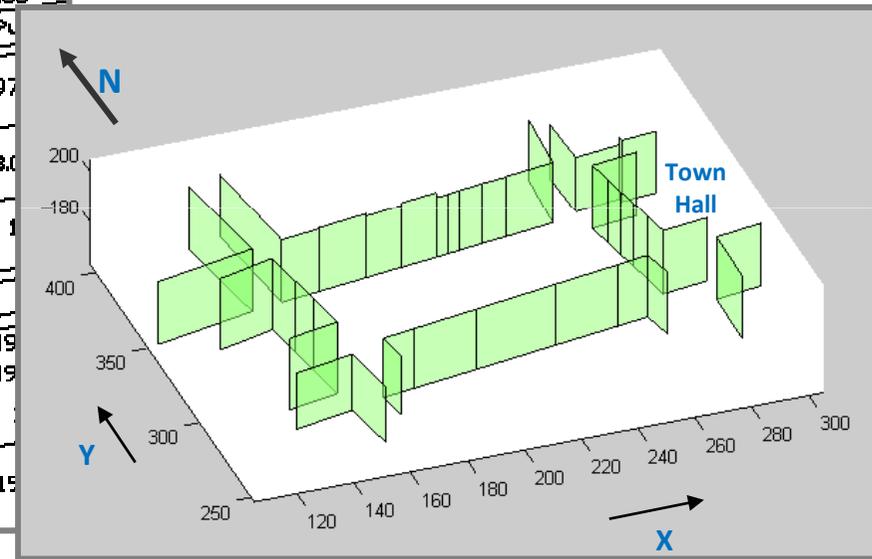
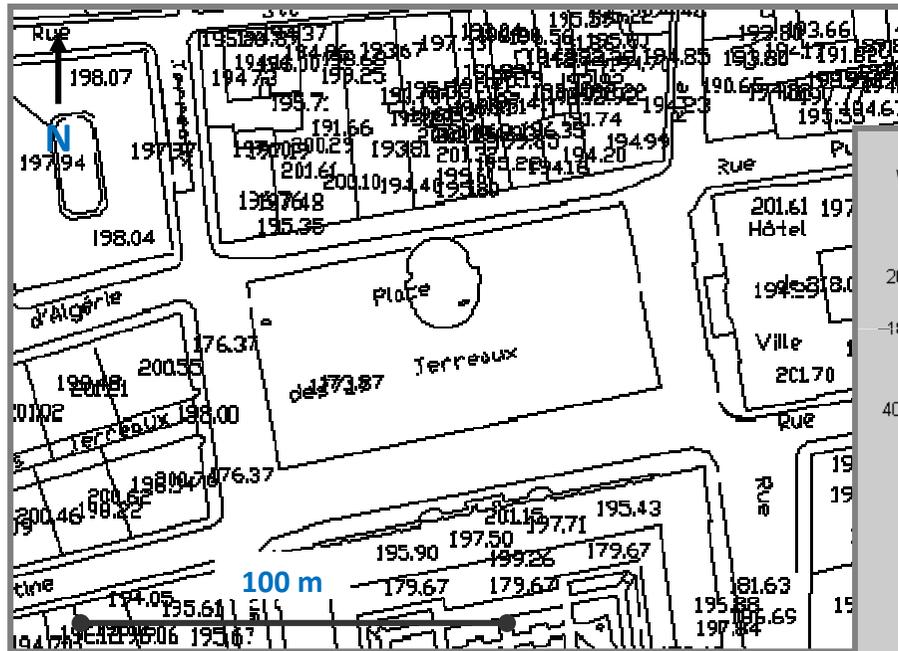
(g)



(h)

Illustration of signature matching using shape contexts and local-neighborhood-graph

Points d'intérêts type SIFT

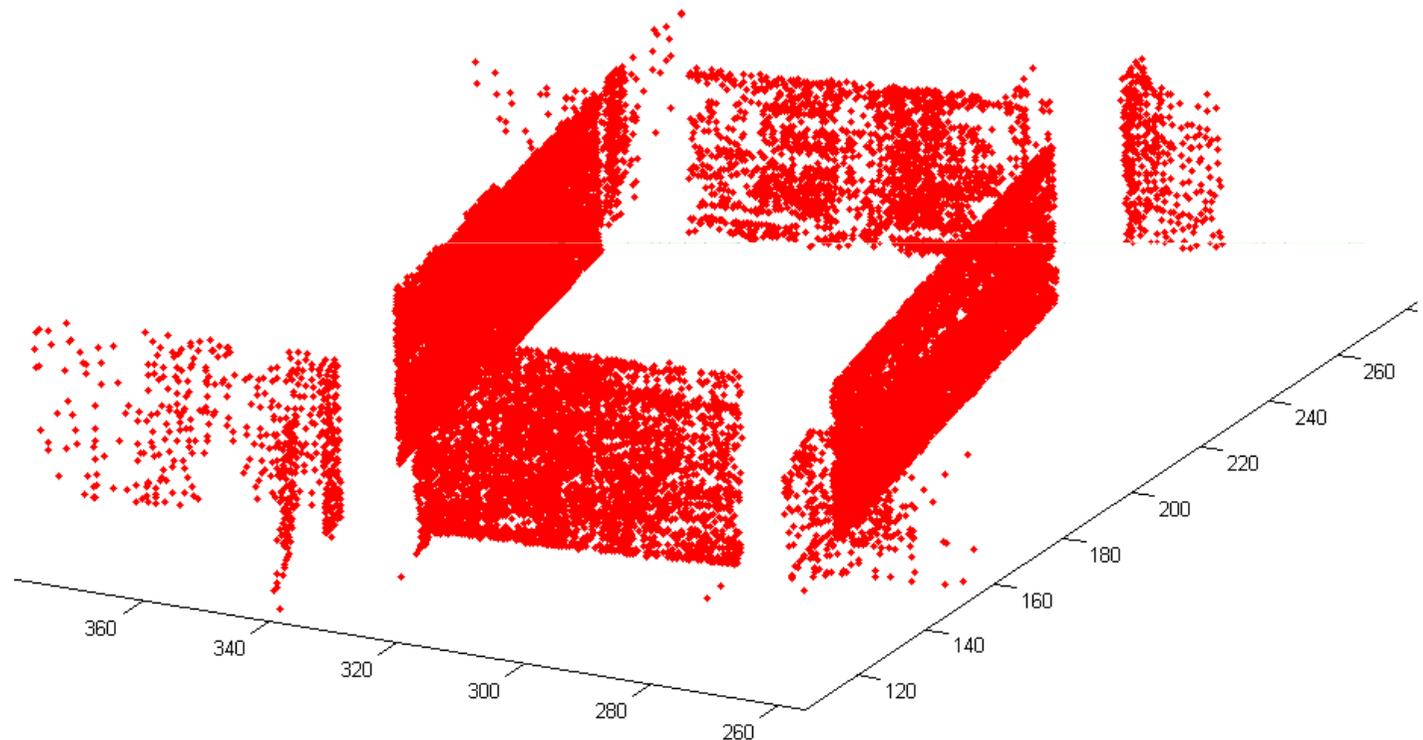


Points d'intérêts type SIFT

36.000 points

Chaque point a :

- une coordonnée 3D
- un descripteur / signature



Points d'intérêts type SIFT

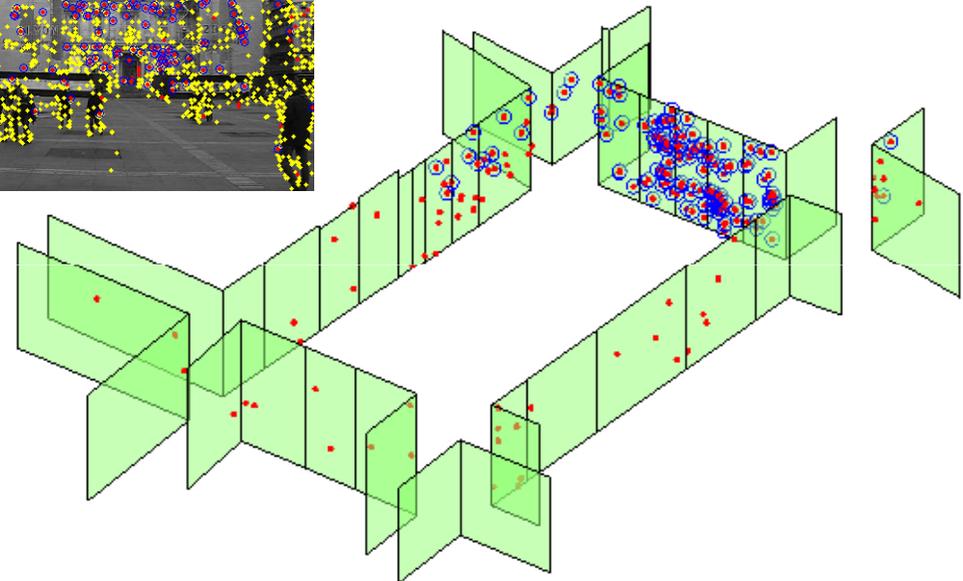
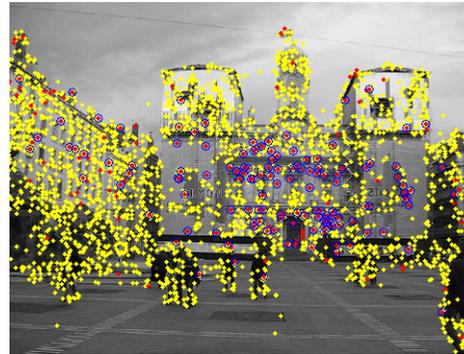
Application sur de la video :

Extraction des points SIFT

Matching avec l'ensemble de référence

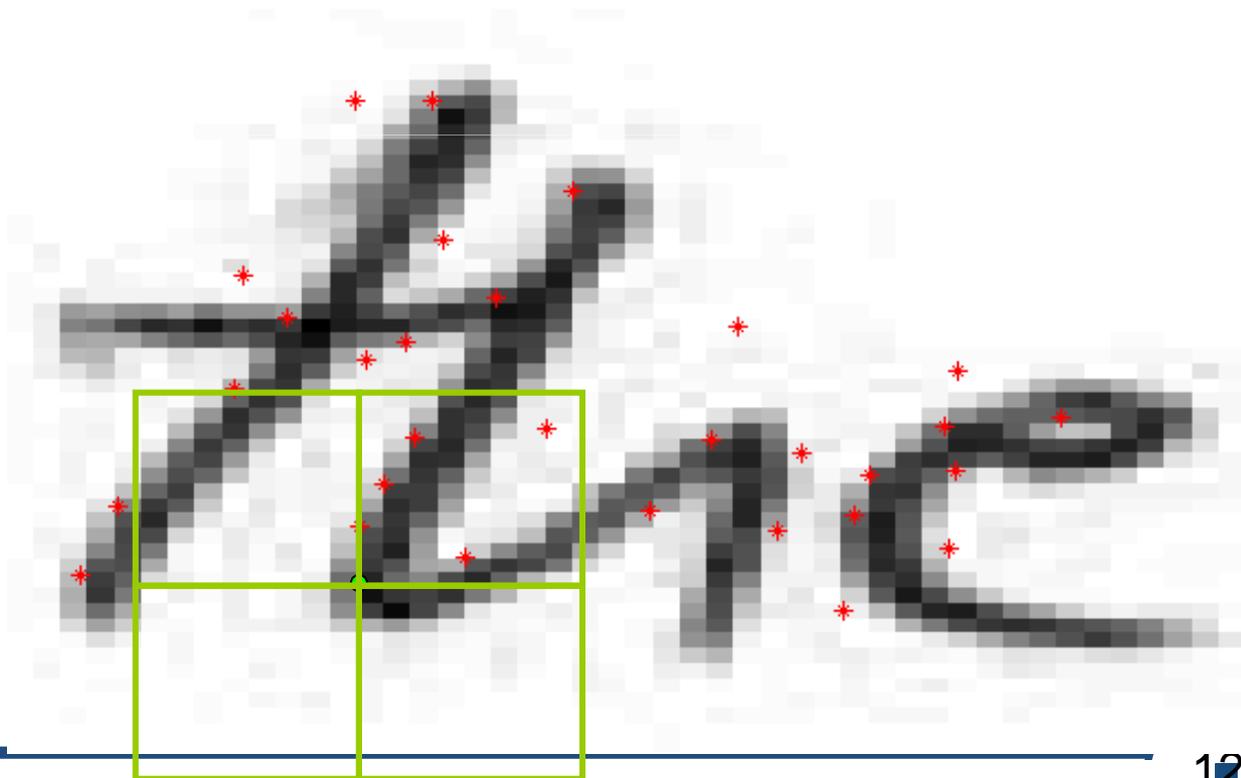
Filtrage des résultats

Appariement avec l'image modèle



Points d'intérêts type SIFT

- Adaptation aux textes :
 - Limitation du nombre d'échelles, de l'invariance en rotation ...
 - Extraction d'information de niveau de gris dans le voisinage des points



Points d'intérêts type SIFT

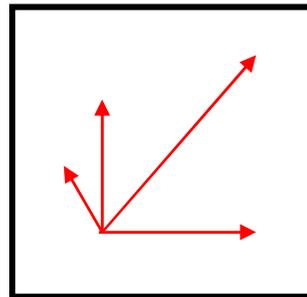
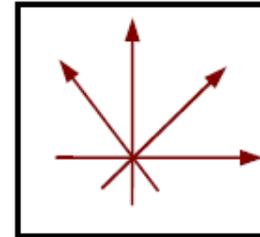
Apart from their formal Admiralty House talks, followed by lunch given by Lady Dorothy Macmillan with Mrs. Kennedy and other guests present, Mr. Kennedy and Mr. Macmillan met three more times yesterday. In PARIS, Mr. Dean Rusk, U.S. Secretary of State, gave a 90-minute briefing on the Vienna talks to the 15-nation Nato council. Some of his listeners said he was

- Limitations :
 - Localisation des points liée aux niveaux de gris et moins à la topologie, géométrie locale
 - Insuffisance du nombre de points d'intérêts

Points d'intérêts type SIFT



Calcul de l'importance
des quatre directions par
auto-corrélation



Adaptation de la signature :

Un vecteur de descripteurs
de dimension 16 pour chaque point d'intérêt.

Points d'intérêts type SIFT

- Résultats intéressants sur la recherche de mots: recherche selon le nombre de points dans l'accumulateur

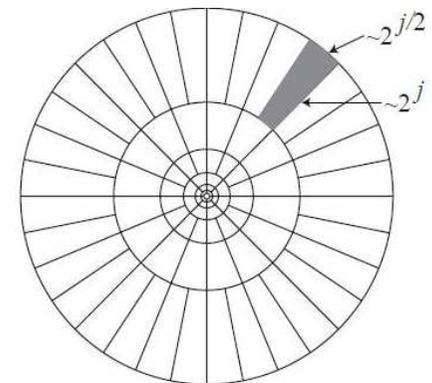
The joint communique on Mr. Kennedy's and Mr. Macmillan's the third talks the first were at Key West, Florida, the second in Washington - said: " Their discussions covered the major problems, both economic and political, and revealed once again the close agreement of the two Governments in pursuing their common purposes. " Occasion was given to review the need for economic co-operation and expansion in the general interests of developed and underdeveloped countries alike. "

Analyse par Curvelets

- **Intérêt pour les informations structurelles des traits**
- **Une caractérisation intégrant à la fois:**
 - **Une information locale** reflétant les structures internes géométriques et anisotropes des traits : courbure et orientation
 - **et une information globale** par accumulation d'infos locales

- **Un choix porté vers un outil de caractérisation mixte: l'analyse par Curvelets**

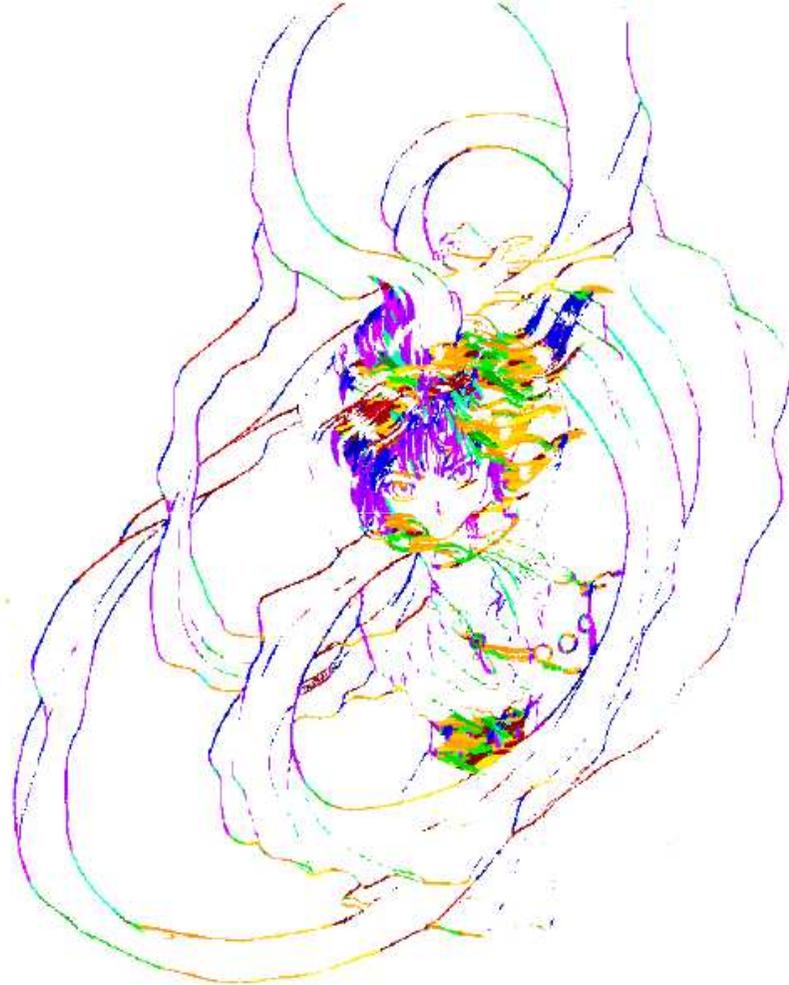
- *Un pavage du plan fréquentiel*
- *Une transformée multi-échelle multi-directionnelle à 3 paramètres: position, échelle et direction*
- *Des fonctions de base de transformées localisées à la fois en espace et en fréquence.*



Analyse par Curvelets



Original image

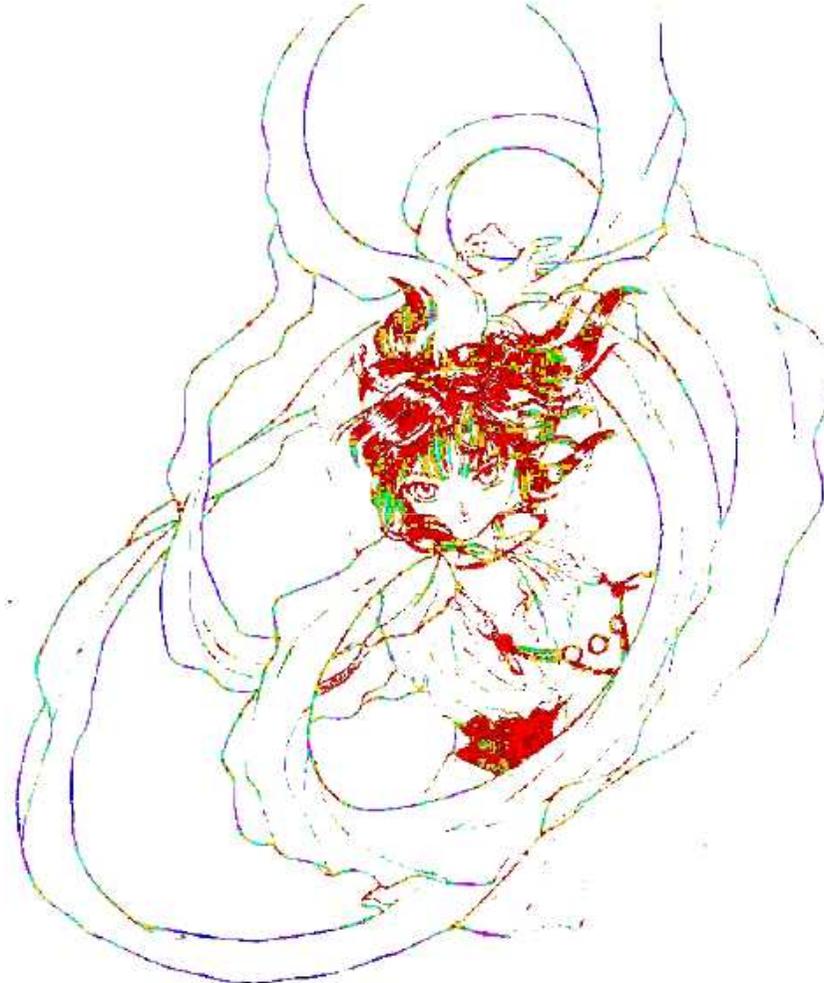


Représentation des orientations

Analyse par Curvelets



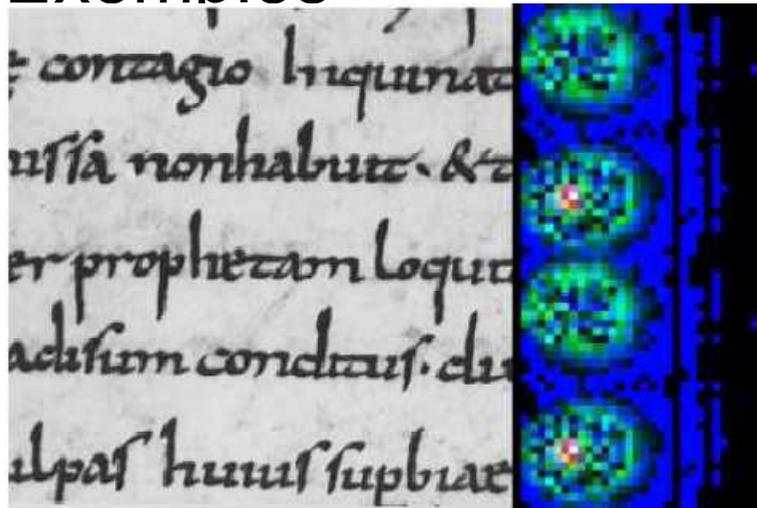
image Originale



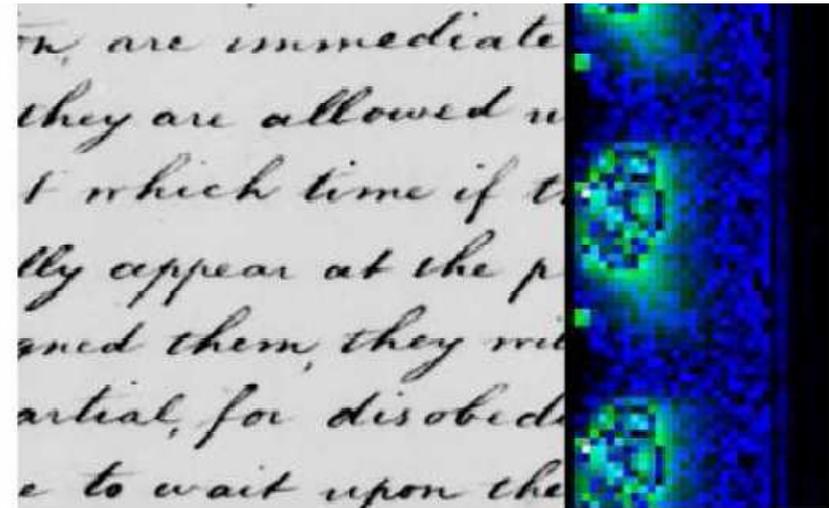
Représentation des courbures

Analyse par Curvelets

- Une signature 2D = Matrice d'accumulation orientation / courbure en tous points
- Exemples



14ème siècle



18ème siècle

Courbure

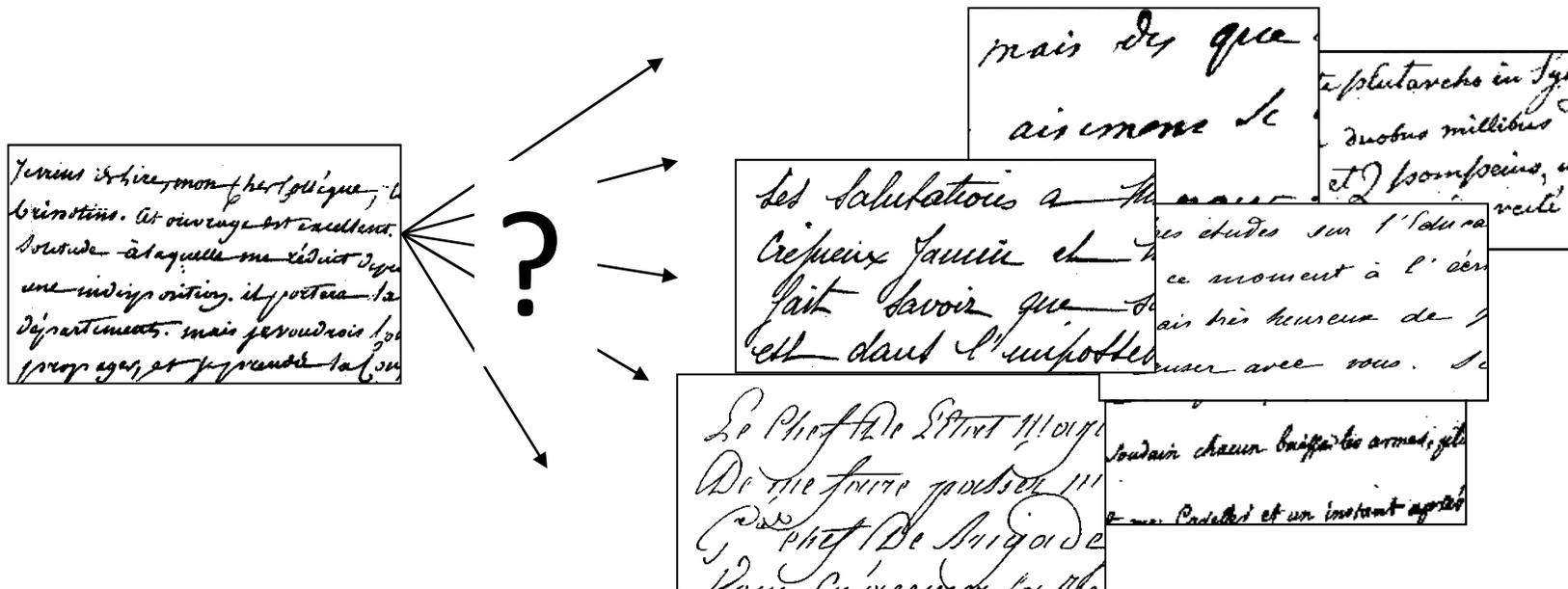
Orientation

⇒ Mise en évidence des propriétés macroscopiques de la géométrie des formes (présence de courbures fortes en des points d'orientations dominantes)

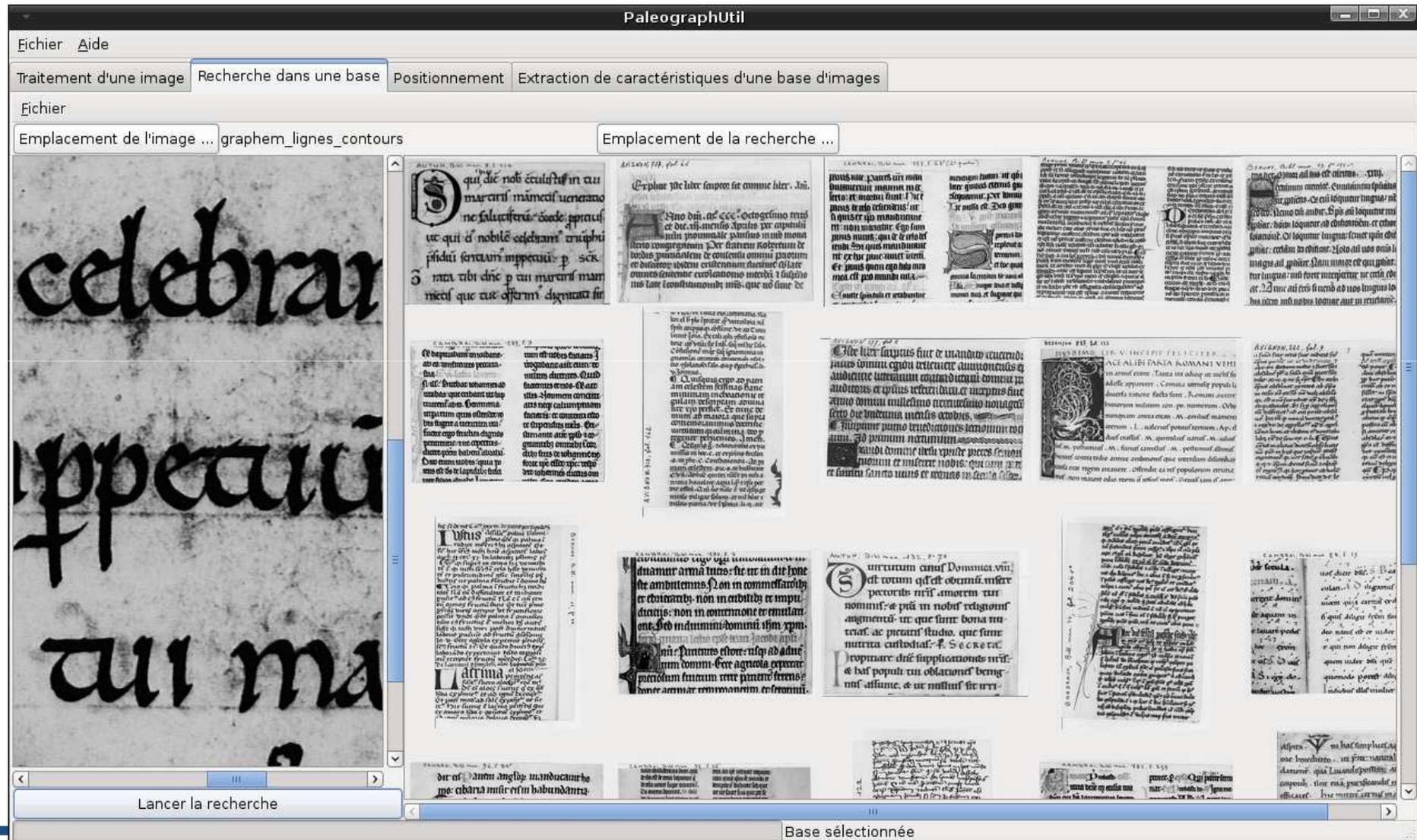
Analyse par Curvelets

- Deux similarités dédiées à deux objectifs: CBIR et groupement perceptif :
 - Similarité générique: corrélation directe entre signatures
 - Similarité « cognitive » (Tversky) non métrique: pondération différente des caractères communs ($Cor_{directe}$) et distinctifs ($Cor_{residus}$)

$$Sim = \gamma \times Cor_{directe} + (1 - \gamma) \times Cor_{residus}$$



Analyse par Curvelets



Analyse par Curvelets

Démo en ligne

The screenshot displays the PaleographUtil software interface. At the top, the title bar reads "PaleographUtil". Below it is a menu bar with "Fichier" and "Aide". A toolbar contains buttons for "Traitement d'une image", "Recherche dans une base", "Positionnement", and "Extraction de caractéristiques d'une base d'images". Below the toolbar, there are two input fields: "Emplacement de l'image ..." containing "Humaniste_petit_sans_binarisation" and "Emplacement de la recherche ...". The main workspace is divided into two panes. The left pane shows a large document page with handwritten text and several horizontal black bars indicating detected features. The right pane shows a grid of smaller document thumbnails. At the bottom, there is a "Lancer la recherche" button and a status bar indicating "Base sélectionnée".

**ET LES APPLICATIONS A
VENIR...**

Vers une gestion électronique des documents complètes et enrichies

E-édition et Chaîne de dématérialisation :

- capteurs intelligents: caméras (tous types de documents), mobilité
- techniques d'indexation (sur les structures, sur les contenus),
- compression des données,
- accessibilité et visualisation des données,
- interactivité numérique,
- pérennité et hébergement des données...



**Archives
Administration
Pédagogie...**

Des besoins spécifiques:

- infrastructure dimensionnée
- interfaces de visualisation interactives
- développement de logiciels ad hoc



**Google books
Digital Libraries
!!!**

Vers une gestion électronique des documents complètes et enrichies

Quelques applications à venir:

Indexation sémantique, navigation personnalisée et visuelle (intuitive) et résumé de contenu

The diagram illustrates two types of navigation in a structured collection:

- Navigation through a structured collection (from title to title for example):** Represented by a red arrow pointing from the title 'LA NATURE' to a specific article.
- Navigation (browsing) in a page (La Nature 1873):** Represented by blue arrows pointing from the title to various sections within the page, such as 'CHRONIQUE', 'CIBIQUES', and 'CIBIQUES'.

The pages shown are from 'LA NATURE' (1873) and contain various articles and sections, including:

- CHRONIQUE:** A section containing news and reports.
- CIBIQUES:** A section containing scientific reports and news.
- LA NATURE:** The main title of the publication.

→ navigation through a structured collection (from title to title for example)
→ navigation (browsing) in a page (*La Nature* 1873)

Techniquement, comment évolue le domaine ?

Vers plus de pluridisciplinarité: nouvelles correspondances entre communautés

- **Images:** passerelles et correspondances d'objectifs, caractérisations locales / globales, recherche objets/mots...
- **Cognition:** modèles cognitifs, interprétation/reconnaissance
- **Perception visuelle:** bas niveau et accès aux données émergentes
- **SHS:** expression des besoins, attentes

Quelles pistes pour aller plus loin ?

- Privilégier les approches sans binarisation ?
- Choix des caractéristiques ?
- Espaces de représentation ?
- Taille des descripteurs ?
- Similarité ?
- Passage à l'échelle ?

....

Résumé

- .Rien de magique dans l'accès aux images, aux contenus
- .L'image seule est une chose, les métadonnées et caractéristiques une autre!
- .Rien ne remplace l'œil, ni le cerveau de l'homme.
Un modèle pour les systèmes actuels