# A DIGITAL APPROACH TO HANDWRITTEN DOCUMENTS
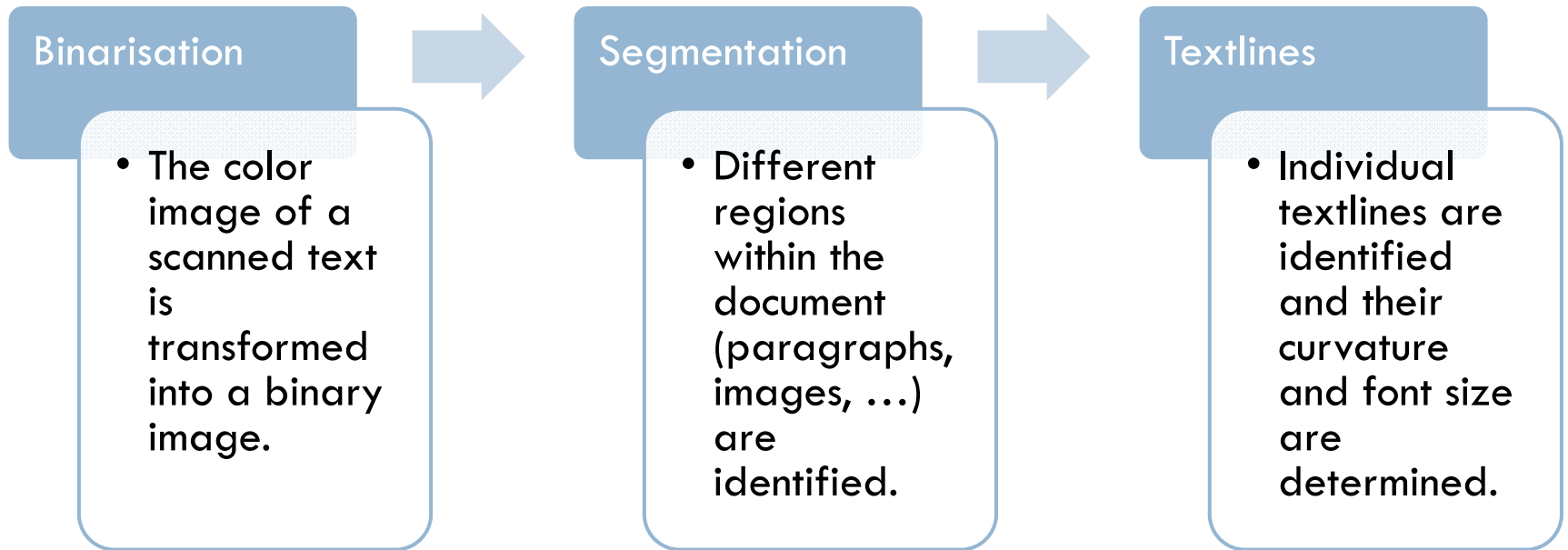
B.I.T. - Bureau Ingénieur Tomasi

# Introduction

- Handwritten documents can for the most part not be read by computers today.

- Our technology such as it has been implemented in the OCR software BIT-Alpha may be a starting point for the development of handwriting analysis tools.

- This digital approach to handwritten documents shall be presented in the following.

# Content capture

- Binarisation
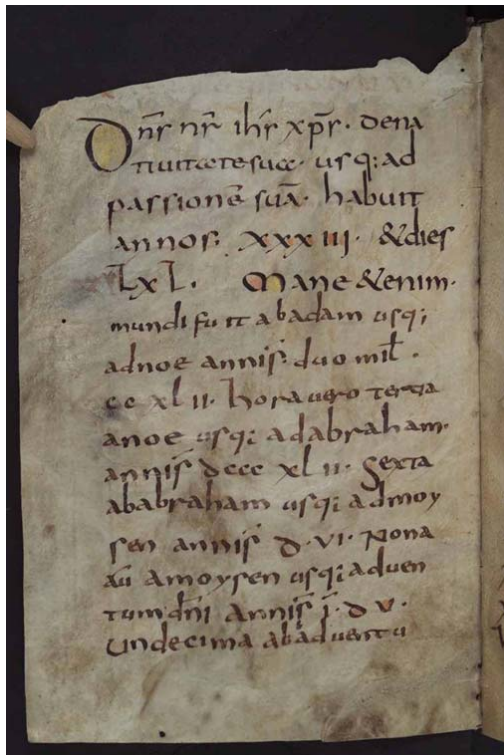
- Segmentation

- Textlines

# Content capture

**Binarisation**

- The color image of a scanned text is transformed into a binary image.

**Segmentation**

- Different regions within the document (paragraphs, images, …) are identified.

**Textlines**

- Individual textlines are identified and their curvature and font size are determined.

# Binarisation

- Archived documents are usually scanned color images (300-600dpi).

- For OCR a binary image of 300-400dpi is required.

- Binarisation is a non-trivial task:
  - Details of the characters are to be preserved.
  - Luminosity and contrast can vary across a single page.
  - The text on the backside of the page may shine through.

# Binarisation

**Scanned image**



**Binary image (mod. Niblack)**

Dnr nr ihr xpr. dena
tiuitatesue· urq; ad
passiona sua· habuit
annos· xxx iij· & dies
Lx L. Mane &enim·
mundi fuit abadam usq;
adnoe anniss duo mil·
cc xl ii· hora uero tertia
anoe usq; adabraham·
anniss dccc xl ii· Sexta
ababraham usqi admoy
sen anniss d vi· Nona
xii amoysen usqiaduen
tum dni anniss i d v·
undecima abaduentu

# Segmentation

- After binarisation different zones within the document are to be identified:
  - Text regions:
    - Titles
    - Paragraphs
  - Graphics regions:
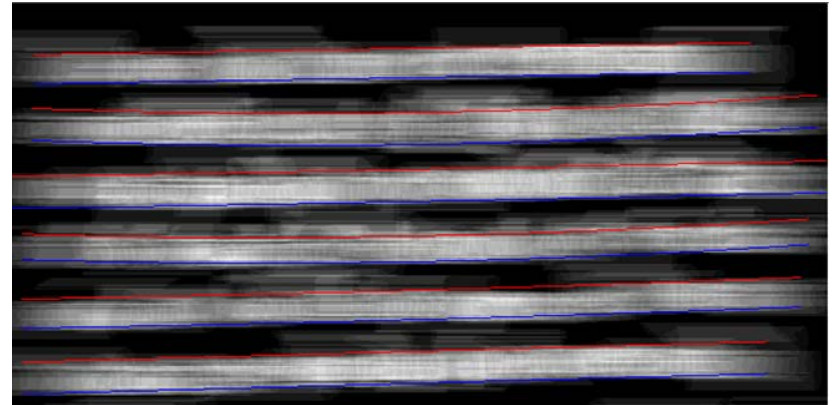    - Images/Figures
    - Lines

# Segmentation

- Segmentation is often combined with a rotation/deskewing of the document in order to obtain horizontal textlines.

- For handwritten documents this can be very difficult: B.I.T. plans to incorporate a bidirectional variance method based on the work of Franck Le Bourgeois (INSA Lyon).

# Textlines

- Capturing textlines is a major difficulty to treating handwritten documents:
    - Textlines may overlap.
    - Characters within a textline can be very close or even touch each other.
    - Textlines may not follow a straight line and be curbed.
- Standard methods for printed documents are not suited to handwritten documents.
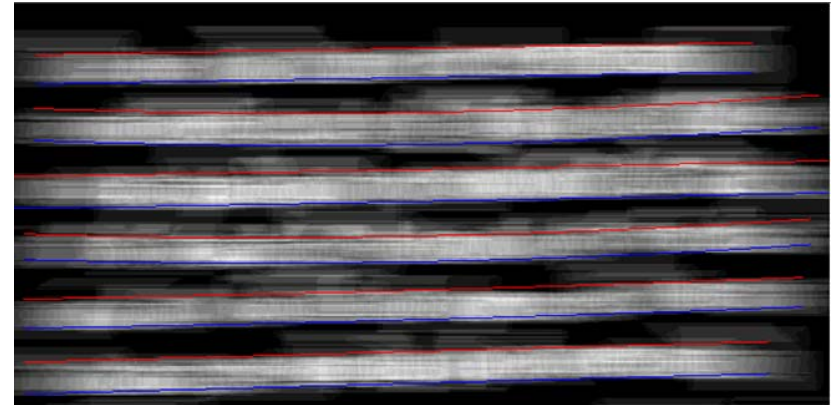
# Textlines

- BIT-Alpha implements a method based on the research of the L.I.R.I.S. laboratory surrounding prof. Emptoz in Lyon: a one-dimensional variation analysis is performed to find zones that constitute textlines.

# Textlines

- After identifying zones of textlines mathematical functions have to be found that fit them.

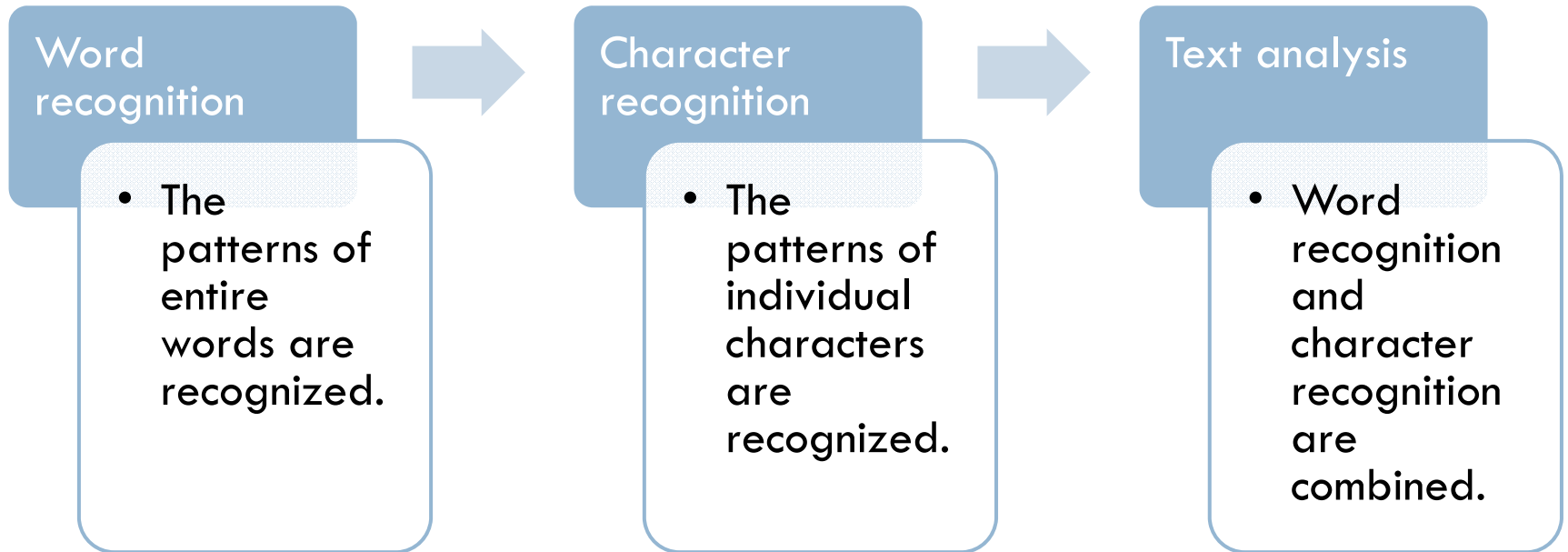- This allows to define a „topline"(red) and a „baseline"(blue) which also provides a measure for the font size.

# Content analysis

- Training and recognition
- Word recognition
- Character recognition
- Text analysis

# Content analysis

**Word recognition**
- The patterns of entire words are recognized.

**Character recognition**
- The patterns of individual characters are recognized.

**Text analysis**
- Word recognition and character recognition are combined.

# Training and recognition

**Training**
- Patterns are trained by the operator.

**Recognition**
- Patterns that are similar to trained ones will be recognized.

# Word recognition

- In handwritten documents individual characters are often joint together and have a high graphical variance

- Identifying entire words can be a natural approach in such cases.

- If not all characters within a word are linked, analyzing the distance between them is required to distinguish gaps between words from gaps between characters.

# Word recognition

- BIT-Alpha has been designed for printed documents of the Renaissance where the distance between words is subject to high variance and not more homogeneous than in handwritten texts.

- The analysis of the distance between characters/words is done on a line-by-line basis in BIT-Alpha.

# Training of words

The pattern of the word „dem" is trained in this example.

# Recognition of words

A new occurrence of a pattern similar to the previously trained one is recognized.

# Character recognition

- In order to obtain a reliable recognition only identifying the patterns of entire words is not enough.

- The isolation of individual characters within word patterns is therefore necessary.

- The approach taken with BIT-Alpha is to recognize sub patterns of individual characters within the patterns of words. We call this process „disjunction".

# Recognition of characters

This screenshot shows the recognition of multiple patterns for the character „a". Note that a blue background indicates a pattern which is identical to a trained pattern, whereas a green background indicates a very similar but not identical pattern.

# Recognition of characters and disjunction

Once the disjunction has been activated other patterns for the character „a" are recognized as sub patterns within words.

## Ligatures

Ligatures can be trained and recognized as a whole.

# Text analysis

- Word recognition proposes strings for each word by comparing the patterns of entire words.
- Character recognition proposes strings for each word by comparing the patterns of constituent characters.
- The correlation of both provides a first clue to the transcription of the text.
- All information concerning the transcription comes essentially from pattern recognition so far.

# Content valorization

- Text transcription

- Text export

- Scribe identification

# Text transcription

- The results from text analysis have to be confronted with linguistic considerations.

- Assuming the language of the text is known the words that have been suggested by text analysis can be matched against a database which contains the words of the language.

- A distance measure for strings is needed to find the best match.

# Text transcription

- The number of edit operations (replacement, insertion, deletion, ... ) required to transform one string into another gives a distance measure („edit distance").

# Text transcription

- The number of edit operations (replacement, insertion, deletion, ... ) required to transform one string into another gives a distance measure („edit distance").

- Errors stemming from OCR-lecture are not necessarily comparable to editing errors (typos). The concept of edit distance has to be extended and adapted for this purpose.

# Text transcription

- In handwritten documents abbreviations have been frequently used.

- We suggest to use Unicode symbols for abbreviations in conformance with the MUFI (Medieval Unicode Font Initiative) recommendations.

- In order to obtain a text that is readily readable by non-specialists these abbreviations have to be expanded.

# Unicode abbreviation

Abbreviations can be trained and recognized as a Unicode symbols.

# Text export

- BIT-Alpha offers a variety of formats for the export of the processed document, including:
  - PDF file containing a color-image of the document in the background and transparent text in the foreground.
  - PDF file containing a binarized image of the document in the background and transparent text in the foreground.
  - PDF file containing only text (opaque).
  - XML file according to the METS/ALTO standard.

## Scanned document

Screenshot of BIT-Alpha showing the scanned document (color image, 300dpi).

## Binarisation and Segmentation

Screenshot of BIT-Alpha showing the binarized document (mod. Niblack) and the results of content capture.

## Text analysis

Screenshot of BIT-Alpha showing the result of text analysis (Unicode symbols are used). The reliability of the recognition is shown by colors (from blue="perfect" to red="doubtful").

## Text transcription

Screenshot of BIT-Alpha showing the transcribed text.

## PDF (binary)

Screenshot of a PDF file with binary image which has been exported by BIT-Alpha. The word „abraham" was searched and has been found.

# PDF (color)

Screenshot of a PDF file with color image which has been exported by BIT-Alpha. The word „principio" was searched and has been found.

## PDF (text only)

Screenshot of a PDF file with only text which has been exported by BIT-Alpha.

# Scribe identification

- In paleography the identity of the person that has written a text (scribe) is of great interest.

- Due to the detailed geometric properties that are delivered BIT-Alpha can be a valuable tool for scribe identification.

- The information provided by BIT-Alpha can be grouped into textline specific and character specific data.

- Such data typically consists of expectation values and standard deviations.

# Scribe identification

- Textline specific data may include (without going into details):
  - Height and thickness of textlines
  - Distance between words
  - Length of words
  - Curvature of textlines
  - Slope of textlines
  - …

# Scribe identification

- Character specific data may include (without going into details):
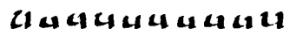  - Size of characters
  - Size of upper/lower-case characters
  - Difference of size between the largest and smallest character within a line
  - Size of accents and i-points
  - Vertical position of accents and i-points
  - …

# Scribe identification

- The pattern recognition capabilities of BIT-Alpha can be used to measure the distance (similarity) between characters written by different scribes.

- Combined with the aforementioned textline and character specific data this suggests the definition of a distance measure for handwritings, which however can only be conceived to be done under the guidance of paleographers.

# Scribe identification

□ The patterns of characters can be exported as bitmaps by BIT-Alpha and will automatically be grouped into subdirectories according to the symbol to which they correspond.

# Conclusion

- The transcription of handwritten documents requires a combination of pattern analysis and linguistic considerations.

- Our approach allows to conceive of a tool that can help with the transcription of texts and scribe identification. The contribution of specialists would be essential, especially to the definition of a distance measure for handwritings.