# eCodicology

Algorithms for the Automatic Tagging of Medieval Manuscripts

# Development of New Technologies for the Automatic Analysis of Medieval Manuscripts

Hannah Busch, Philipp Vanscheidt (University of Trier)

Swati Chandna (Karlsruhe Institute of Technology)

Celia Krause (Technical University of Darmstadt)

# Project and Motivation

# Project and Motivation



**What is eCodicology?**
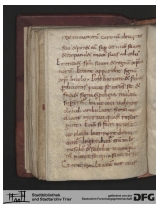
A BMBF funded joint research project of
- Technical University Darmstadt
- Trier Center for Digital Humanities
- Karlsruhe Institute for Technology

# Project and Motivation

## What is eCodicology?

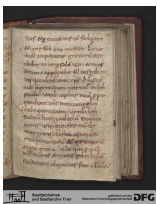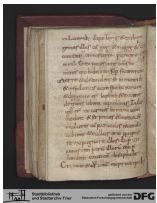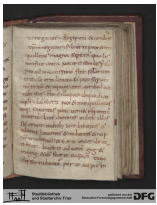A BMBF funded joint research project of

- Technical University Darmstadt
- Trier Center for Digital Humanities
- Karlsruhe Institute for Technology

## What is eCodicology for?

- Automatic identification of macro and micro structural layout elements
- Quantitative Codicology: Statistical analysis of reproducible features
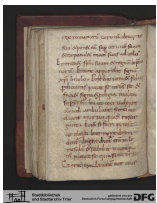- To identify hidden relationships
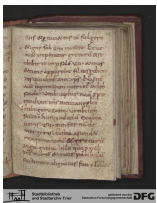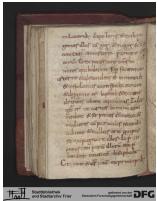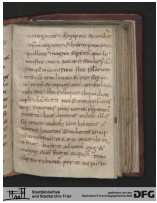
# How to Publish a Library?

Virtual Scriptorium St. Matthias

# How to Publish a Library?

Virtual Scriptorium St. Matthias



Scope:    440 codices online
Period:    8th to 16th century
Contents: Liturgica, Bible, patristics,
           mysticism, philosophy,
           poetry, law, charters, etc.

# How to Publish a Library?

Virtual Scriptorium St. Matthias

TextGrid



Scope: 440 codices online
Period: 8th to 16th century
Contents: Liturgica, Bible, patristics,
mysticism, philosophy,
poetry, law, charters, etc.

# How to Publish a Library?

## Virtual Scriptorium St. Matthias

## TextGrid

Scope:    440 codices online
Period:   8th to 16th century
Contents: Liturgica, Bible, patristics,
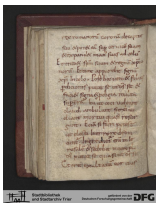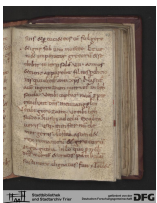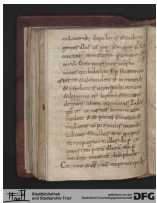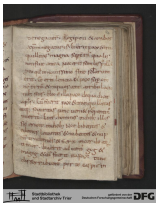          mysticism, philosophy,
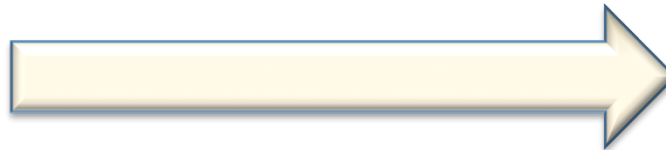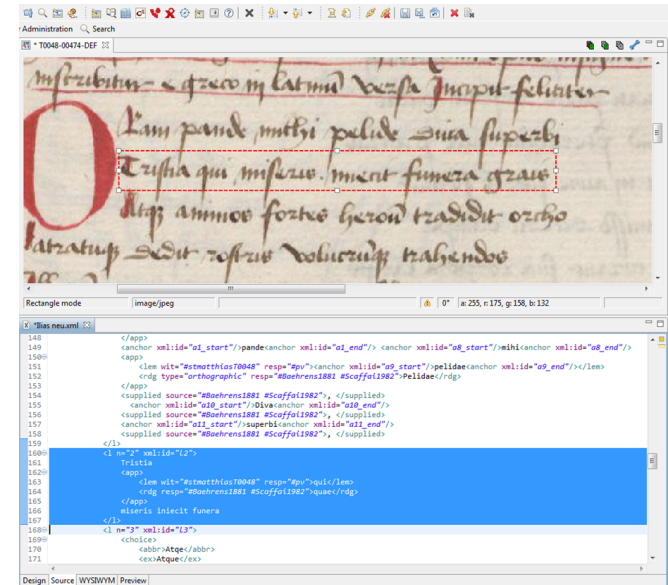          poetry, law, charters, etc.

KoLibRi

Mass Ingest

# How to Publish a Library?



**Virtual Scriptorium St. Matthias**

Scope: 440 codices online
Period: 8th to 16th century
Contents: Liturgica, Bible, patristics, mysticism, philosophy, poetry, law, charters, etc.

KoLibRi

Mass Ingest

**TextGrid**

- PID Handling by GWDG
- Exploration of material by human researchers, eg. Critical editions

# Feature Extraction

◆ **Feature:** Individual measurable property of phenomenon observed



page size,    written size,    pictorial space

# How to Extract Layout Features?

Color Calibration → Preprocessing → Segmentation and Feature Extraction → Data Storage

# Color Calibration

- Adjusting different color spaces to similar color space according to a standard color chart



Before Calibration

Scanner 1

Before Calibration

Scanner 2

After Calibration

Scanner 1

After Calibration

Scanner 2

# **Preprocessing**

◆ To increase accuracy of the digital data for feature extraction

    ◆ Spatial Calibration - Adjusting different resolutions to same resolution

    ◆ Filtering - Removing low frequency background noise

    ◆ Scaling and Duplication – Scaling images to different resolutions



300 dpi        400 dpi Different Resolutions      Spatial Calibrated Image

# Segmentation

◆ Divide image into regions with similar properties

◆ Determine region of interest

# Data Storage

◆ Extracted features are stored in XML file



```
<xml>
....
<width>250
</width>
<unit>mm
</unit>
<height>550
</height>
<unit>mm
</unit>
…</xml>
```

# **Feature Extraction….In a nutshell**

◆ Automatic layout analysis of manuscript images is challenging due to different color spaces, heavy noise etc.

◆ Suitable preprocessing and segmentation methodology are very essential steps for feature extraction

◆ Accuracy and convergence rate of such techniques must be significantly high in order to ensure the success of subsequent steps

# Metadata Management

- Customization of TEI, using TEI ODD – One Document Does It All

- Preferable schema language for output: RelaxNG + Schematron

- Special focus on measurements, written and pictorial spaces, e.g. marginalia, pictorial elements

```
<xml>
....
<width>x
</width>
<height>y
</height>
...
</xml>
```

```
<tei: TEI>
....
<tei:width>x
</tei: width>
<tei:height>y
</tei:height>
...
</tei: TEI>
```

# Metadata Management

TEI XML Example for Manuscript Description:

```xml
<tei:TEI>
    <tei:teiHeader>
    …
        <tei:msDesc>
        …
            <tei:supportDesc>

                <tei:extent>
                    <tei:measure type="leaves">44</tei:measure>
                    <tei:measure type="format">8°</tei:measure>
                </tei:extent>
                …
            </tei:supportDesc>
            <tei:layoutDesc>
                <tei:layout>
                    <tei:dimensions type="written" corresp="#written1">
                        <tei:height quantity="250" unit="mm" min="249" max="251" confidence="0.8">250mm</tei:height>
                        <tei:width quantity="100" unit="mm" min="98" max="101" confidence="0.77">100mm</tei:width>
                    </tei:dimensions>
                    <tei:dimensions type="image" facs="#image1"/>
                </tei:layout>
                <tei:layout columns="1" writtenLines="26 30"/>
            </tei:layoutDesc>
            </tei:objectDesc>
            …
            <tei:additions>
                <tei:note type="gloss" place="gutter"/>
            </tei:additions>
            …
        </tei:msDesc>
        …
    </tei:teiHeader>
    <tei:facsimile>
        <tei:surface type="leaf" ulx="0" uly="0" lrx="100" lry="250">
            <tei:graphic url="0099-00005.jpg"/>
            <tei:zone type="image" ulx="20" uly="20" lrx="70" lry="70" xml:id="image1"/>
            <tei:zone  type="written" ulx="20" uly="20" lrx="90" lry="90" xml:id="written1"/>
        </tei:surface>
    </tei:facsimile>
    <tei:body>…</tei:body>
</tei:TEI>
```
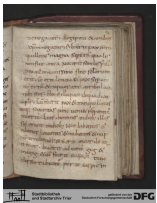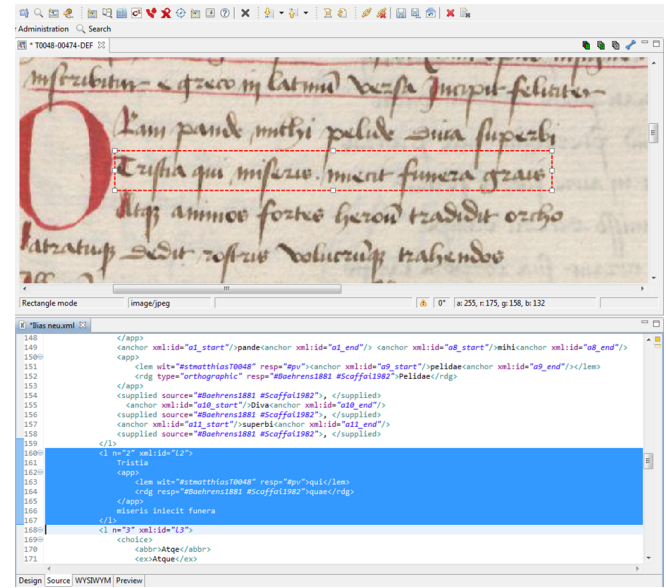
# How to analyze the results?

# Quantitative Codicology

# Current Case Study    (100 codices)

## Numerical variables          (absolute values)

A)    page dimensions          (height  x  width)

B)    number of columns

⎤ 20 randomly selected pages

C)    number of folios

D)    number of blank pages within a codex

## Categorial variables

E)    codex format          (2°; 4°; 8°)

F)    dating

G)    text genre

H)    writing material          (parchment; paper; both)

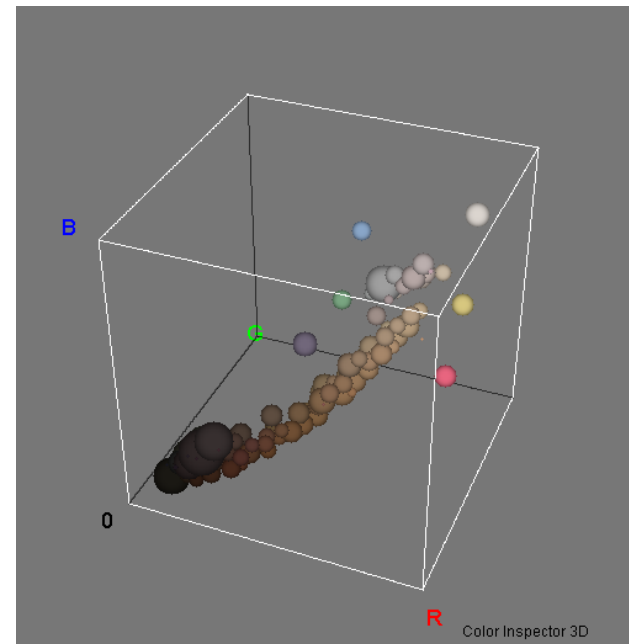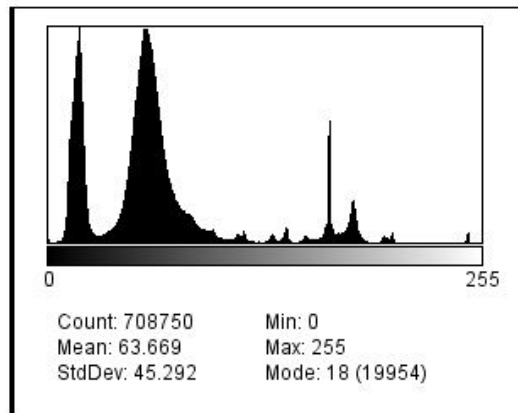| Handschrift | Textgattung / Titel | Beschreibstoff | Datierung | Format | Kennung der Seite | Spaltenzahl | Seitengröße (mm) | |
|---|---|---|---|---|---|---|---|---|
| Signatur, z.B. S Hs 144 | | | ohne AD | | z.B. S Hs 144_017 | Nummer | Höhe | Breite |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 223,632 | 131,411 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,897 | 134,253 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,161 | 131,678 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,897 | 133,517 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 223,264 | 134,989 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 3 | 223,264 | 135,356 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 222,161 | 133,149 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 221,793 | 135,356 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 222,161 | 134,989 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 224 | 136,092 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 223,632 | 134,989 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 224,368 | 137,931 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 2 | 223,632 | 135,724 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,161 | 137,563 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 219,954 | 136,828 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 219,954 | 136,092 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,161 | 134,253 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 221,057 | 137,195 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 222,161 | 136,828 |
| T Hs 1128/2053 | onentia /De aequivocis / I | Papier/Pergament | 14.Jh.; 15.Jh.; 16.Jh. | 8° | rier.de/stmatthias/T1128/T1 | 1 | 221,425 | 135,724 |

Handschrift (37) / Handschrift (38) / Handschrift (39) / Handschrift (40) / Handschrift (41) / **Handschrift (42)** / Har



0 — 255

Count: 708750  Min: 0
Mean: 63.669  Max: 255
StdDev: 45.292  Mode: 18 (19954)



Color Inspector 3D

# Conclusion

Our goals are …

to build correlations between different parameters within groups of **codices with common features** (for ex. dating / text genre / writing material)

to make comparisons between groups of **codices with different similarities** (for ex. codices made of parchment vs. codices made of paper or handwritten codices vs. incunables)

The eCodicology project is associated to