# TexT – Text Extractor Tool for Handwritten Document Transcription and Annotations

Anders Hast^, Per Cullhed* & Ekta Vats^

^Department of Information Technology,

Division of Visual Information and Interaction
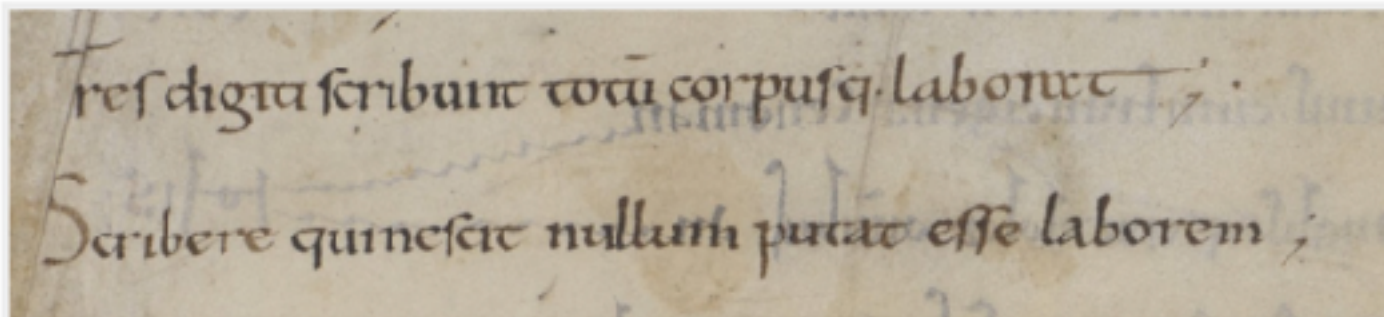
*Uppsala University Library

# Outline

- Transcription is a **tedious** task that can be made **simple** by handwritten text recognition techniques, such as **word spotting**

- It will be shown how transcription can be done in a **new way**, but also how researchers can get **easy access**, using similar techniques, to manuscripts that are otherwise not searchable

- Semi-Automatic transcription

  - In arbitrary word order instead of line by line, word by word

  - (Problem solving through a visualisation of the problem)

# Copying books is an old trade

© LONDON, BRITISH LIBRARY, ROYAL 6 A.VI, FOL. 109R

*Tres digiti scribunt totum corpusque laborat. Scribere qui nescit nullum putat esse laborem.*

*[Three fingers write and the whole body labours. He who does not know how to write thinks it is no work.]*



## Marginalized

*Notes in manuscripts and colophons made by medieval scribes and copyists*

- New parchment, bad ink; I say nothing more.

- I am very cold.

- That's a hard page and a weary work to read it.

- Let the reader's voice honor the writer's pen.

- This page has not been written very slowly.

- The parchment is hairy.

- The ink is thin.

- Thank God, it will soon be dark.

- Oh, my hand.

- Now I've written the whole thing; for Christ's sake give me a drink.

- Writing is excessive drudgery. It crooks your back, it dims your sight, it twists your stomach and your sides.

- St. Patrick of Armagh, deliver me from writing.

- While I wrote I froze, and what I could not write by the beams of the sun I finished by candlelight.

- As the harbor is welcome to the sailor, so is the last line to the scribe.

- This is sad! O little book! A day will come in truth when someone over your page will say, "The hand that wrote it is no more."
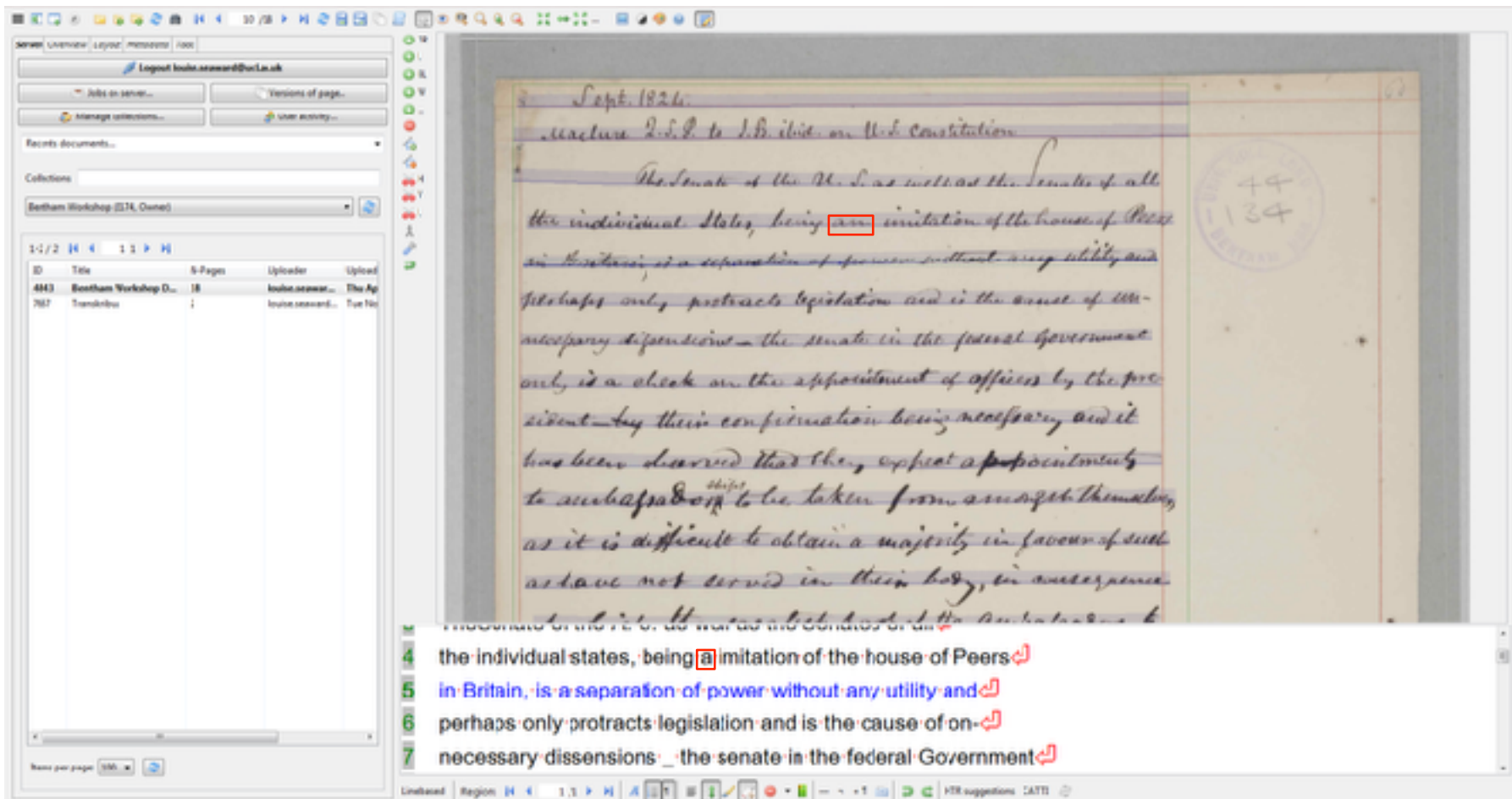
- The monks tried to brighten up their day by scribbling in the margins
- If copying the text was boring, transcribing it will be likewise boring…

# Transkribus

https://transkribus.eu/Transkribus/

**Even modern transcribers makes mistakes…**

# Transkribus

- Now, this sounds interesting, but it starts to get really exciting if you consider that
    - **once you have properly transcribed e.g. 100 images** you may inform us and we will train an HTR engine from the Computational Intelligence Technology Lab (CITlab) of the University of Rostock on your documents and
    - you will be able to transcribe further pages of your documents with the **support of automatically produced handwritten text**.
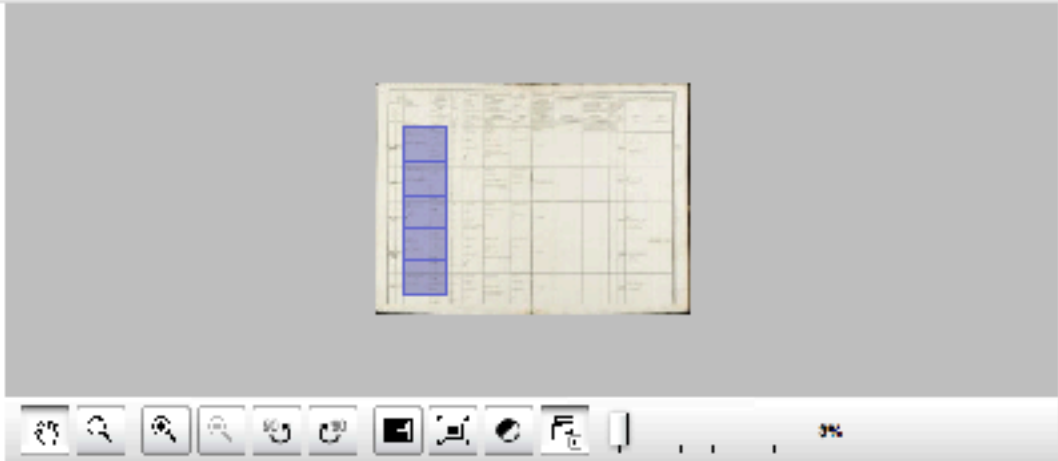
# Veele handen

Is a Dutch initiative where volunteers can sign
in and start transcribing in a simple manner

# T-pen

From the university of St. Louis in the United States is much used for the transcriptions of medieval manuscripts.

# The Zooniverse transcription project

The Oxford based Zooniverse is a general crowdsourcing platform used for many purposes, not only transcription. It has one million users-mostly from English speaking countries.

## Operation War Diary

Every transcriber knows exactly where to start contributing as this is obvious from the GUI

# The Purpose of Transcription

- Our goal with this project is to make our collections:
    - more accessible and useful to curators, researchers, and anyone with a curious spirit
- Because computers have a hard time understanding handwriting, many of our collections still hold many secrets and hidden knowledge inside their pages
    - With your help, we can bring that knowledge to life.
- Searchability & Readability - Usefulness over Perfection

**38% COMPLETE**

100 Total Pages
27 Contributing Members

**START CONTRIBUTING TODAY.**

**LEO BAEKELAND DIARY VOLUME 29, 1920**

How did a scientist split time between social calls, "automobiling," and laboratory work at the beginning of the twentieth century? Help us transcribe Leo Baekeland's diaries to learn more about his daily activities and scientific work.

# TexT-project - Uppsala

**(1) Transcribe**

Volunteer(s) work together to transcribe document. When finished, can request review.
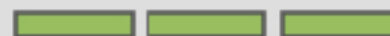
**(2) Review**

A new volunteer reviews transcription. Can approve or send back to (1) for more edits.
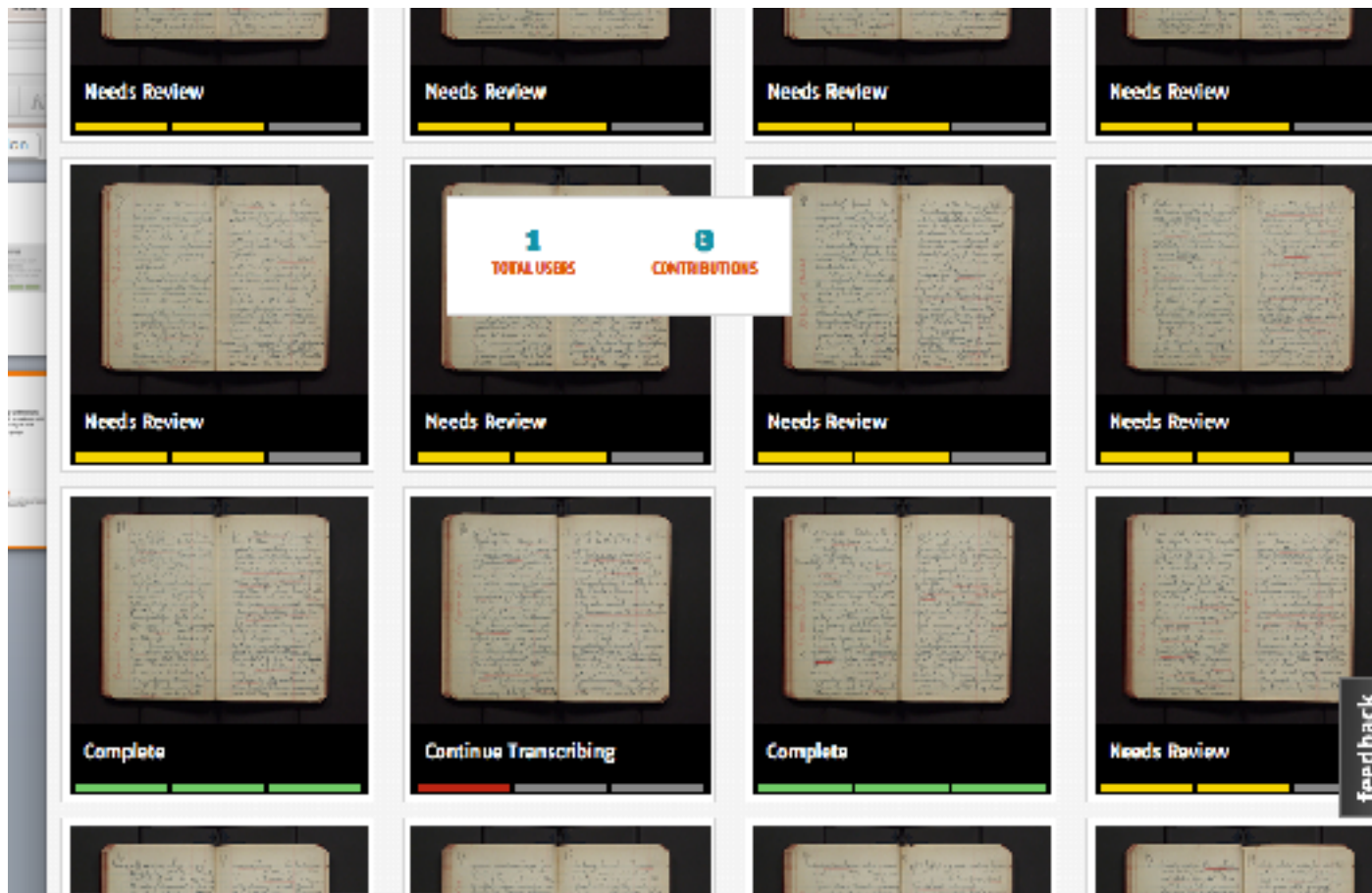
**(3) Approve**

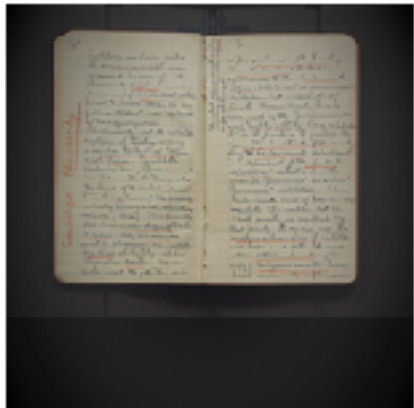Smithsonian staff approves transcription or sends back to (1) for more edits

HTR

Machine learning

- Why do people love Sudoku?!?
- Is it really a game?
  - Or rather about problem solving?
  - And the satisfaction of reaching a goal?
- You learn NOTHING?!
- From crosswords you might at least learn something…

# Can Transcription be "fun"?

- Can we make it a game?
  - How and why? Should we even??
- Or is it rather about problem solving?
- Can you learn something while doing it?
- This is a challenging task!

- **Gamification** is *the application of game-design elements and game principles in non-game contexts*

- Probably no first person shooter game…

  - not even medieval fighting games… or?!
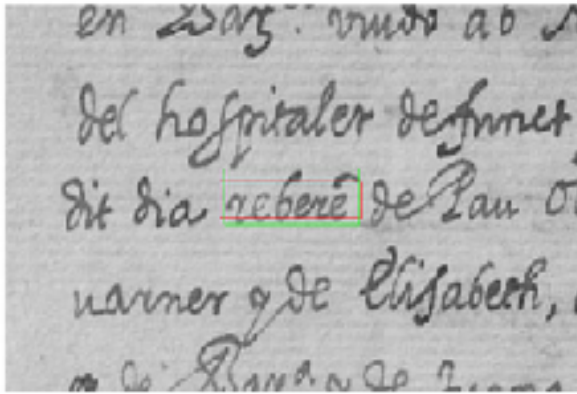
- What about problem solving?

  - If Sudoku can be fun..
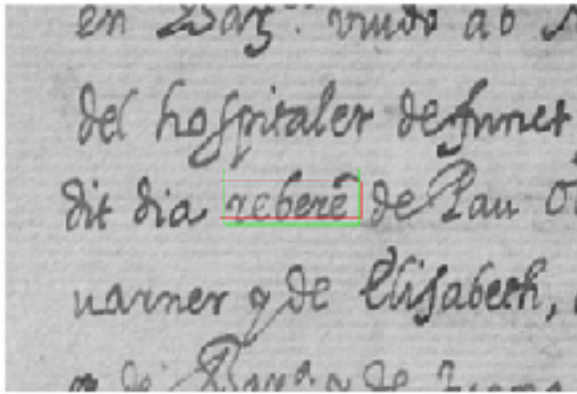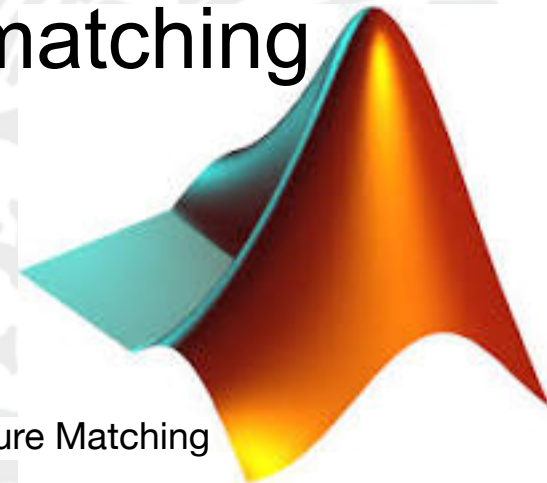
- The user marks up a word
- …and the HTR engine finds it in the text
- Can be used to make large text document collections searchable
- Ex: Monk
  - L. Schomaker, Rijksuniversiteit Groningen
  - Search & annotation tools for handwritten manuscripts
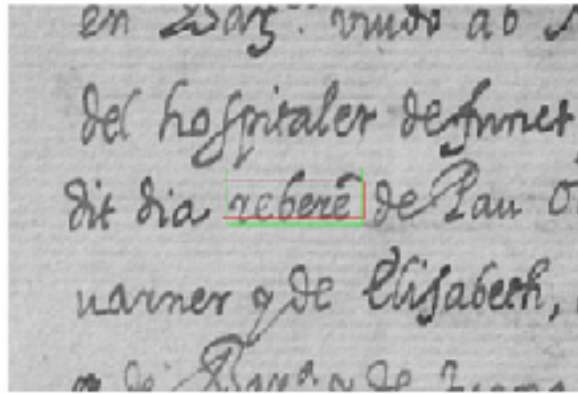- Monk learns by user annotation

**Word Spotting**



- We have developed an HTR engine for Word Spotting (WS)
- Segmentation Free
  - Based on stitching principles from Computer Vision
  - i.e pattern matching
- Learning free

A Segmentation-free Handwritten Word Spotting Approach by Relaxed Feature Matching
Anders Hast and Alicia Fornés
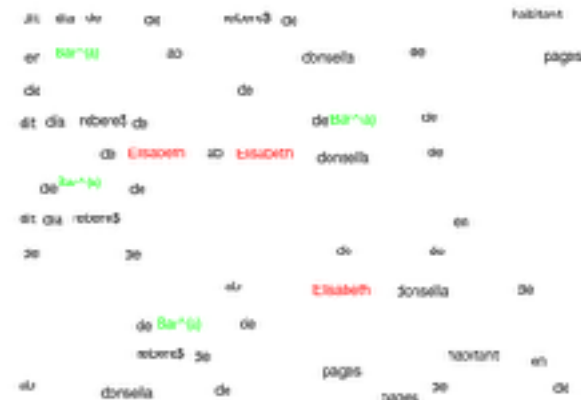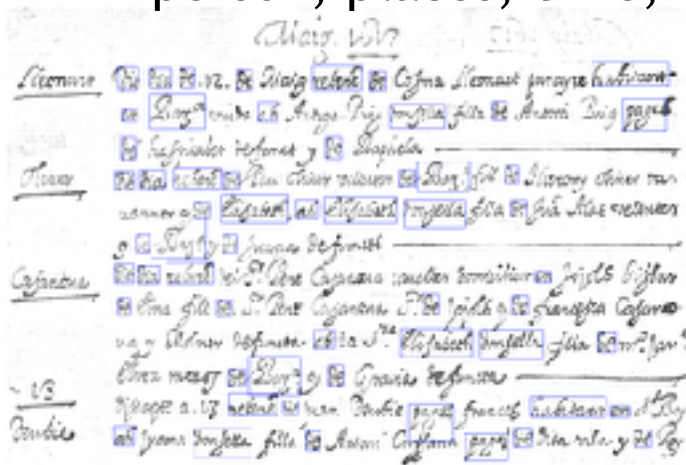DAS2016, International Workshop on Document Analysis Systems, 2016.

- We plan to use the WS approach for transcription as well!

- The user marks the red rectangle
- The system finds the "perfect" bounding box
  - Principle: "all the word and nothing but the word"!

# Semi Automatic Transcription

- User should mark a word and transcribe it ONCE

- Then the HTR-Engine should find all occurrences

- Annotation possibilities!

  - person, places, time, etc…

TexT - Text Extractor Tool for Handwritten Document Transcription and Annotation
Anders Hast, Per Cullhed and Ekta Vats
14:th Italian Research Conference on Digital Libraries, IRCDL 2017

# Problems

- Only whole words are found
- Note how "d'Audenarde" and "a Audenarde" are contracted to one word
  - This is usually a small problem
- A bigger problem are the split words
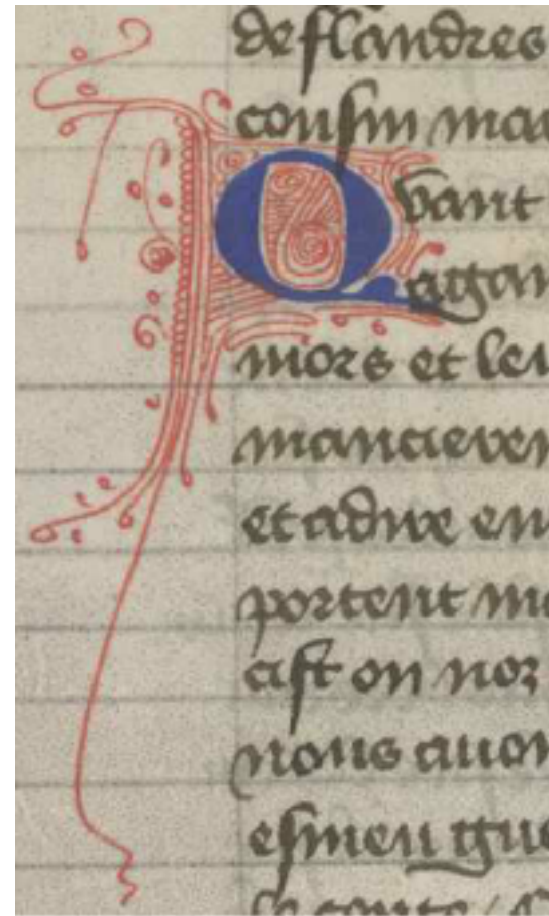- Solvable by letter spotting?

# Problems with Initials
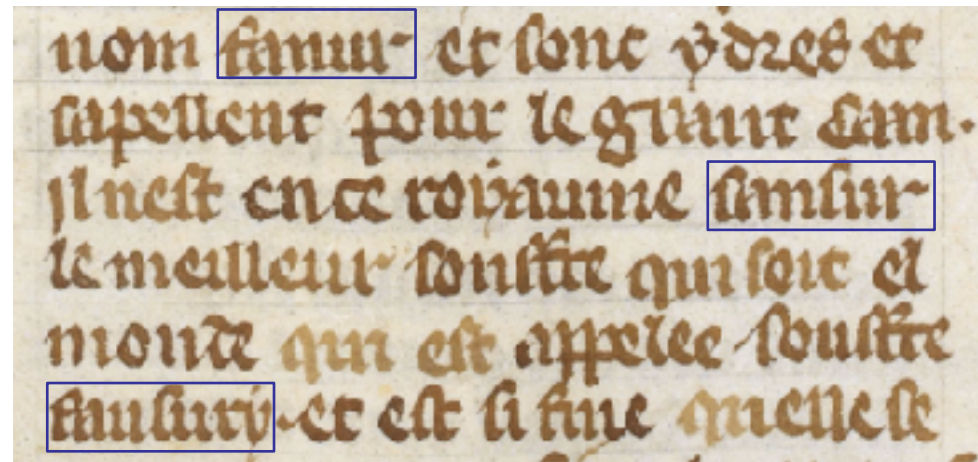
**Close to searched word (Audenarde)**

**Or if it is the word we search (Quant)**

# An Example of problems with Orthography

- A section of Marco Polo's account dealing with Fansur

- Here spelled three different ways:

    - <fanur>

    - <sansur>

    - <fausur(y)>



British Library Royal MS 19 D I, f.123r

https://medium.com/@siwaratrikalpa/barus-sumatra-in-a-medieval-egyptian-text-edcfb7be9517

# Abbreviations is not a new thing

## The most common texting abbreviations, Italian-style

- m = **mi** (*I, me*)
- t = **ti** (*you*)
- xke = **perché** (why, because)
- cmq = **comunque** (*anyway*)
- bc = **baci** (*kisses*)
- midi = mi dispiace (*I'm sorry*)
- pfv = per favore (*please*)
- d = **da** (from, since, of)
- grz = **grazie** (*thanks*)
- tn = **tanto** (a lot, much, long time)
- k = **chi** (who, what)
- c6 = **Ci sei?** (Are you there?)
- qls = **qualcosa** (*something*)
- + = **più** (*more*)
- risp = **rispondi** (*answer*)
- nn = **non** (*no, not*)
- prox = **prossima** (*next*)
- gg = **giorno** (*day*)
- tvb = ti voglio bene (*I love you*)
- ta = **ti amo** (I love you)

## Common abbreviations, Monk-style



http://www.dummies.com/languages/italian/texting-and-chatting-in-italian/
https://kuscholarworks.ku.edu/bitstream/handle/1808/1821/47cappelli.pdf

# Actually it's not a new thing at all

"M[arcus] Agrippa L[ucii] f[ilius] co[n]s[ul] tertium fecit"

Marcus Agrippa, son of Lucius, made [this building] when consul for the third time

# Text and Writing Variations

# Query Expansion

- Search for a word

- Repeat the search for each and every word found

- Old thing… but we do not have a feature for each word to use in order to improve there result

- Use the Visualisation tool: Beatrice

# 'Human in the loop'

# Results

# Why?

- Use the tool to find a cut between inliers and outliers
- Examine further:
    - the worst 'correct' and
    - the best 'incorrect'
    - in doubt: look at the sentence where the word occurs
- Let the system LEARN from this
- The advantage is that it is MUCH easier than looking through the 50 pages to check…
- Non found words will be "left over" after complete transcription

# Conclusion

- WS can be used to make non transcribed collections searchable
- It can also be an efficient tool for making transcriptions easier
- Use visualisations to **solve** the problem:
  - correct vs. incorrect words
- Still many interesting challenges!