uwe springmann
❧ scientia cum humanitate ❧
digital humanist

**Wittscholarship 2014**

# ❧Latin Natural Language Processing
# ❧OCR-Nachkorrektur-Tool PoCoTo

2014-06-06

# Latin NLP?

- huge heritage: largest body of historical literary sources
- Latin publications dominate print production until about 1750
- many titles have never been reprinted
- either key or barrier to cultural heritage of the western world
- has been left out of the EU IMPACT (**Imp**roving **Ac**cess to **T**exts, 2008-12) project despite its importance

# Foresti 1486, Supplementum Chronicarum

- color is gone, but "printing space economy" still holds (abbreviations)
- was heavily used by Schedel's Weltchronik (Latinist project @LMU)

# Foresti 1486 (very modest training)

frugales: eo r tyrãnidisexptes ex agro oĩa ad vlctũ necessaria preter oleũ abũdãtissime babeãt. Multi.m/

0001/010007.bin.png

thebeoum legionis in ea glorisso martyrio coronati sunt: quorum corpora ẽt nunc p̃eipuo venerãtur honore/

0001/010008.bin.png

uGenuit hec ip̃a cioitas maximũ ipsius ep̃m: virum vtique scĩa et sanctitate p̃eipuũm: qui in suis codicibus: vt./

0001/010009.bin.png

li.ɔ. loco suo diceĩ: multa cdidit.Jn huius agro ad cinisij montis radices secusia fuit ciuitas: quã Federi

0001/01000a.bin.png

cus barbarossa per mõte3 illũ trãsiens dolo captã diruit: et nũ φ̃ in eodẽ stata coaluit: cuius ep̃atus sedes:

0001/01000b.bin.png

postφ̃ et ip̃i ciues episcopu3 suũ trucidarunt: in alterius dictionẽ cessit.

0001/01000c.bin.png

@Qporedia cisalpine gallie baud magna vrba a R0. populo õ gallos infestos: vt plinio placet: sbyllinis li/

## Adam von Bodenstein, 1577

**Kreüter**

ner erſcheinüg/ vnſerer teütſcher
zaun oder hagwurtzel/ gar nicht/
welche der mehrtertheil balbierer
für rechte Ariſtolochiam rotun=
dam einſamlend. Dioſc. Diſer
wurtzel etwas mit wein myrrhen
vnd pfeffer getruncken/ reiniget
die weiber von vberfliſzigem vn=
rath der můter/ treibt auſz die an
d geburt vñ weiber menſes. Ein
ſalb gemacht vonn diſer wurtzen
zeitloſen vñ anagallide zeücht vfz
fpreiffel/ dȯrnvñ geſchiferte bein.
Hiemitt beſchlies ich mein
rede diſer zeit von den zwelff zei=
chen kreütteren/ begåren menck=
lich welle mirs im beſten aufnem
men als dañ ichs gethan/ hab fye
weitleüffiger beſchribẽ wellen/ fo
find yetziger kürtze viel vrſach_en/
voraufz dieweil ich groffen koften
angewendet in fůchung der kreü
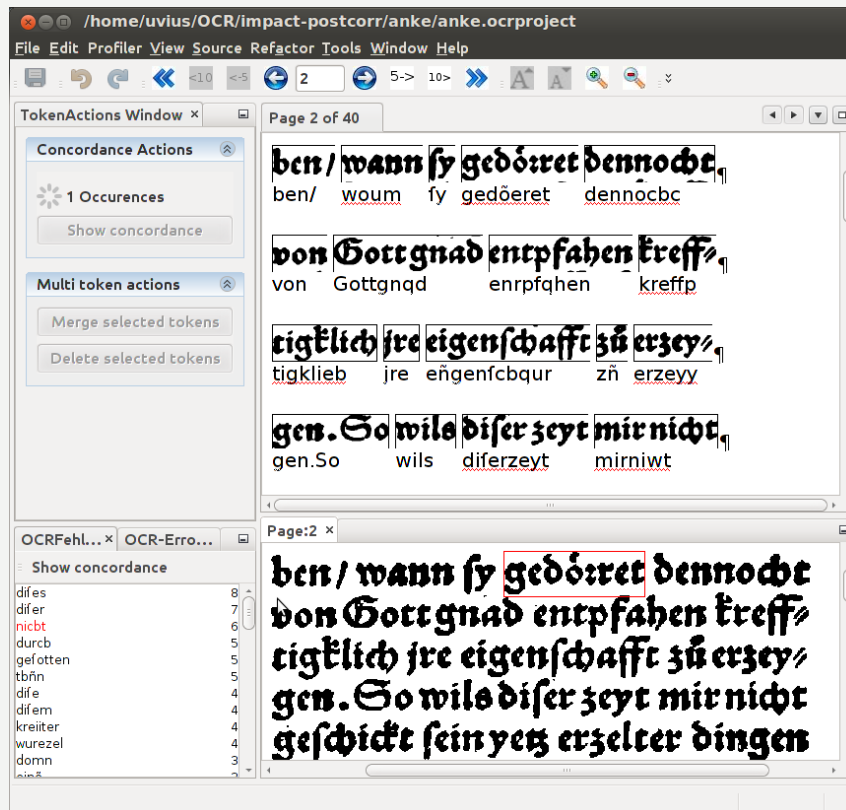ter aufz eignem willen vñ beüt_tel/
ncn=

# PostCorrectionTool PoCoTo

- locally installable Java package for postcorrection
- developed at CIS as part of the EU IMPACT project
- word synopsis: image + OCR with interactive correction
- error profiling: calculate statistical error model based on
  a) historical spelling (not an error)
  b) proper OCR errors
  and propose most probable correction candidate
- batch correction: rank errors according to frequency & error pattern
  and enable quick correction decision in concordance view

- try it out:

http://www.digitisation.eu/tools/browse/ocr-post-correction-and-enrichment/post-correction-tool/
https://github.com/thorstenv/PoCoTo

# PoCoTo (developed at CIS)



- error frequency
- word synopsis
  (tesseract hocr output)
- page context

# PoCoTo (developed at CIS)

error pattern concordance with batch correction

# Thank you for your attention!