

Historical legislation

CASE STUDY



Annemieke Romein
Sara Veldhoen



KB } nationale
bibliotheek

About us

Annemieke Romein

- Researcher
- KNAW Huygens ING
- Researcher-in-Residence at KB in 2019



Sara Veldhoen

- Research Software Engineer
- KB, national library of the Netherlands



Content

- Hypothesis
- Sources in Numbers
- From scans to computer-readable text
- Categorisation
- Training Annif

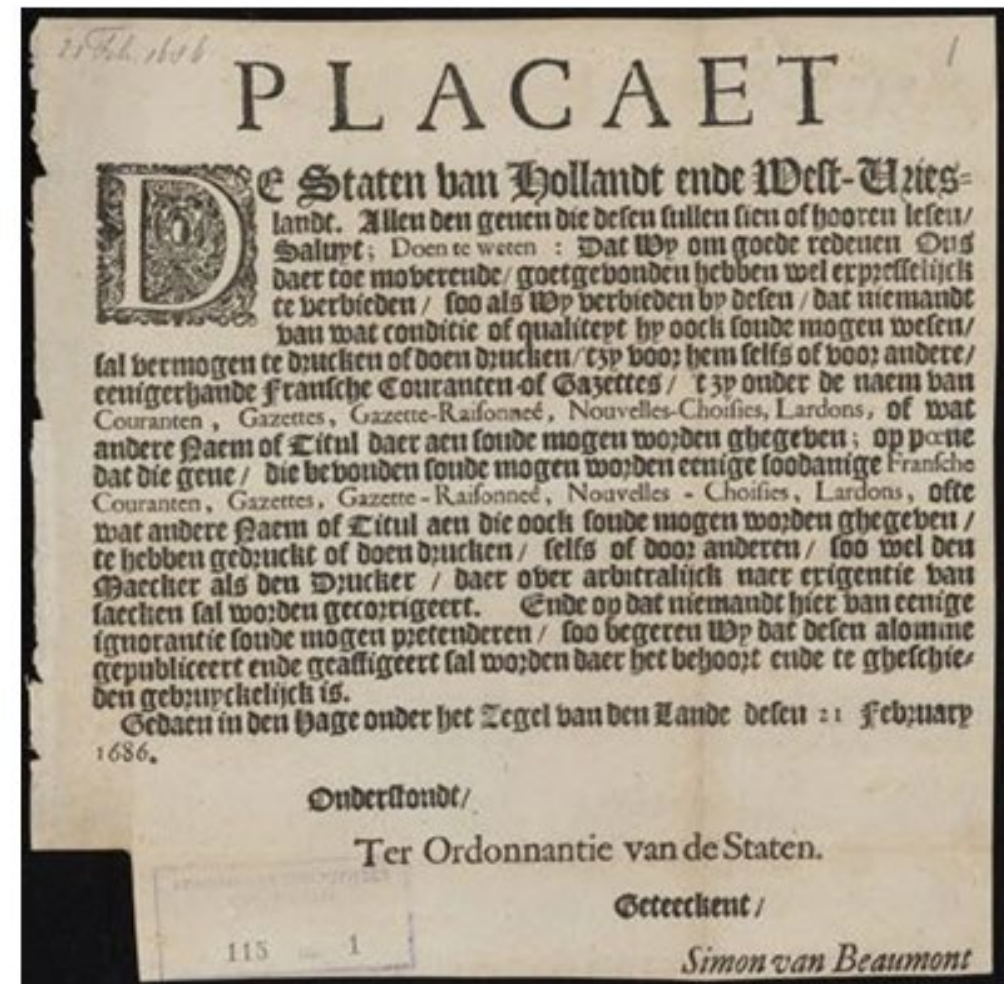




Paalhuis,
Amsterdam
(ca. 1660)

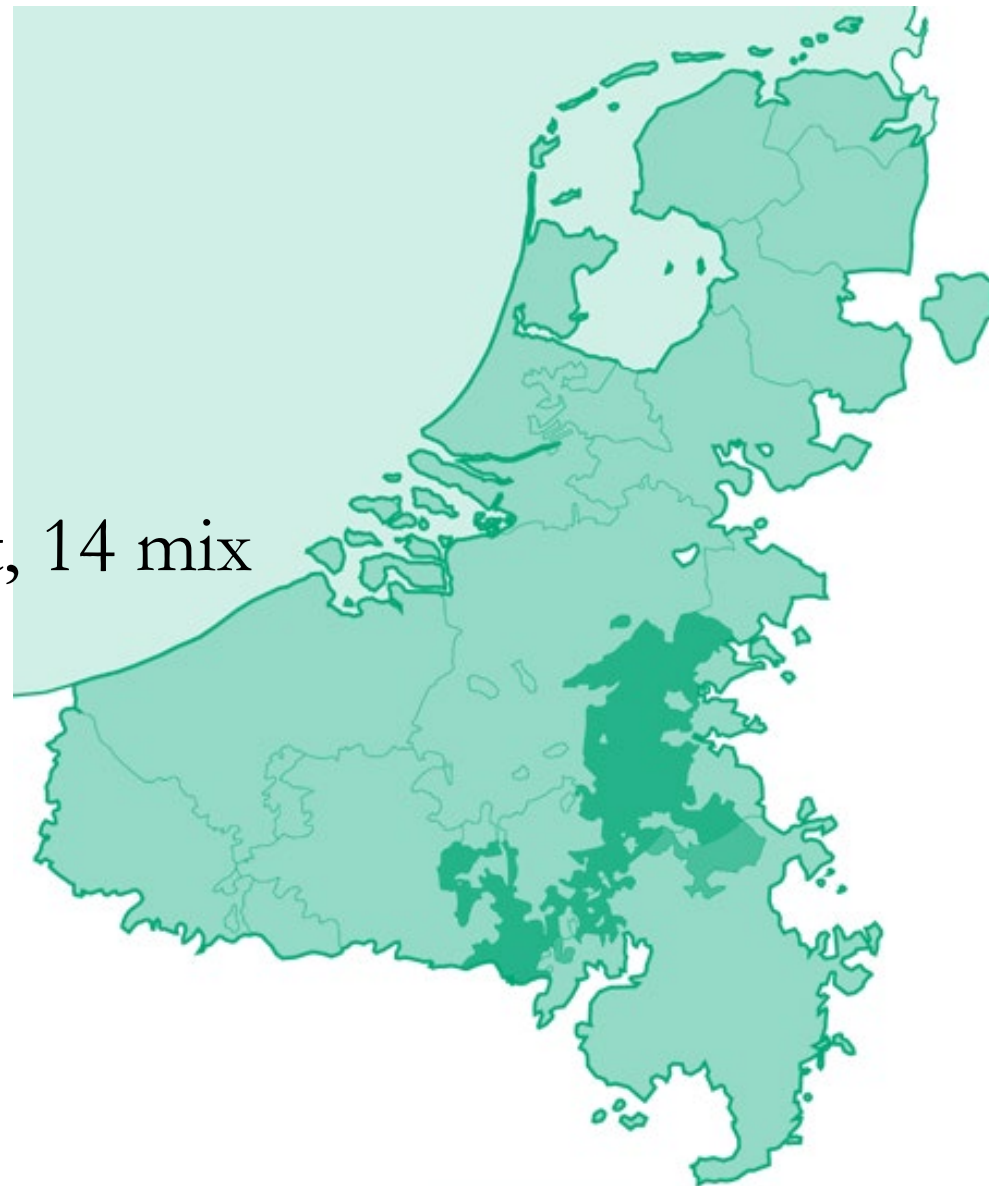


From volumes to individual laws



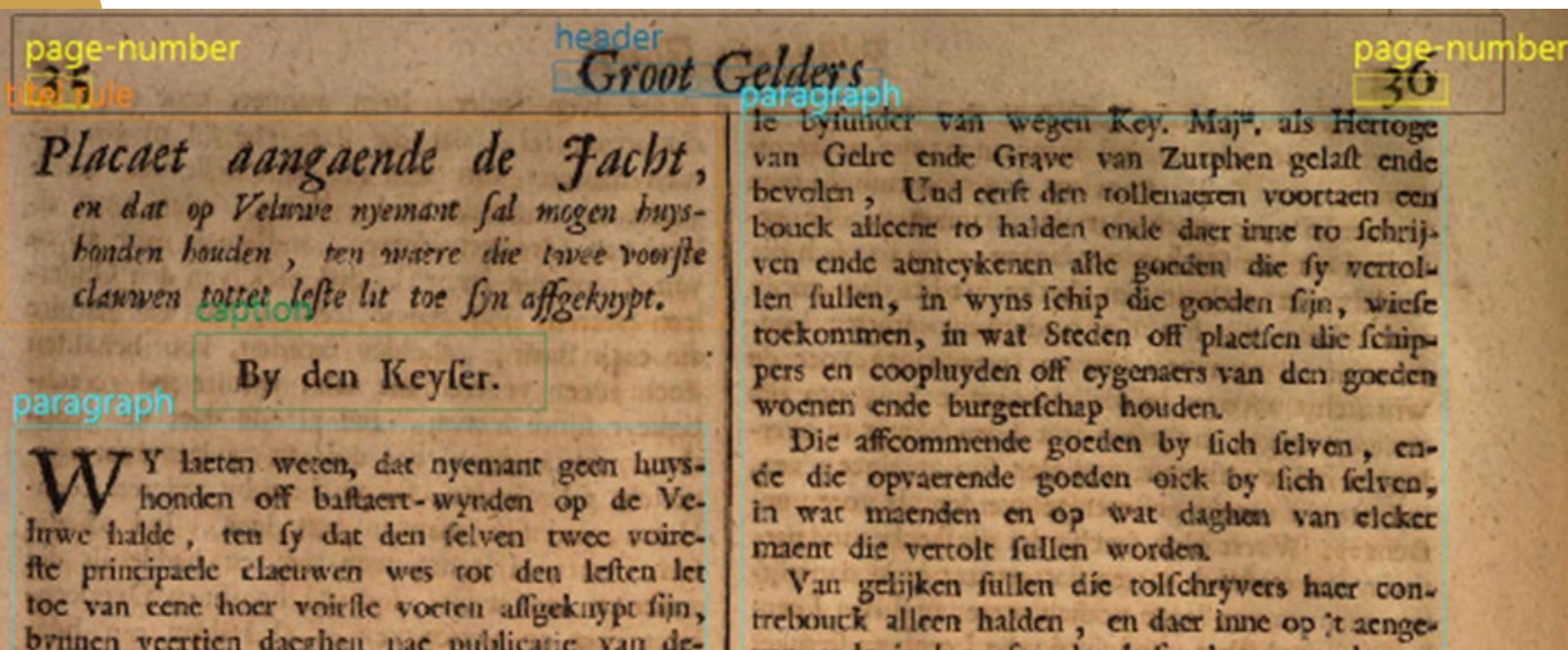
Sources in numbers

Sources #	108
Font	88 roman, 20 gothic
Language	67 Dutch, 26 FR, 1 Lat, 14 mix
Pages #	75.000
Characters#	550 mln
Publication date	17th and 18th century
Region	Nearly all provinces



Teaching the computer to read and to segment texts

Transkribus



D Staten

van Stadt Gronin- gen ende Ommelan- den: DOEN TE WETEN. Alhoewel wy verhoopt hadden dat door het redres en correctie van de Lijste van verpachtinge by ons in het laest verlo- pen Jaer 1622. gedaen ende alomme gepubli-

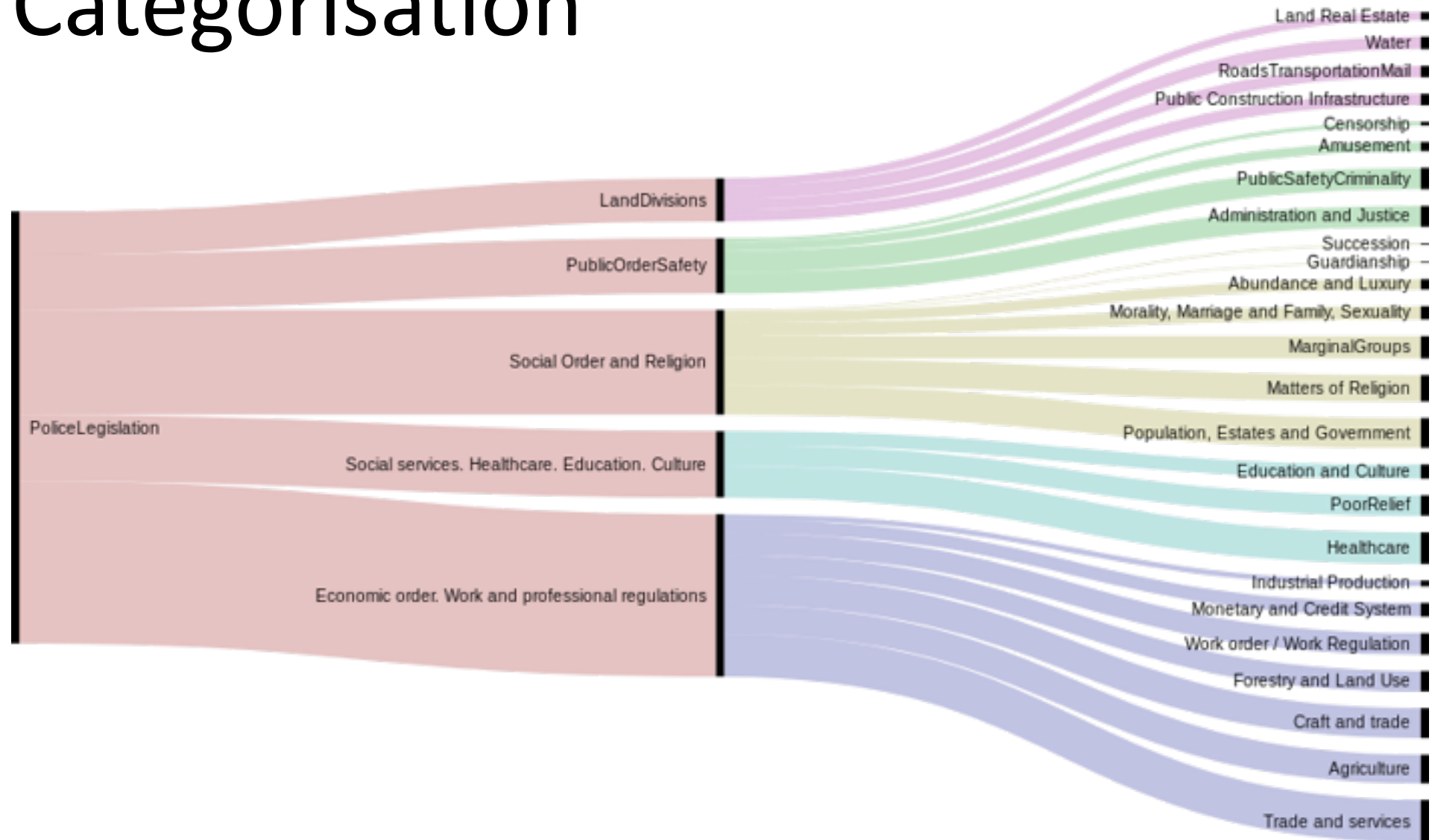
ceert / eenighsins gheremedieert ende belet sulden zijn geweest / de veelvuldige frauden ende dieberpen / Item listige ende schadelijcke practiquen / compositionen en an- dere abuyfen / streckende tot groten aff-breuck van de Generale Middelen / en merckelijcke prejuditie van alle goede getrouwe Ingesetenen ende oprechte Liefhebbers van 't Vaderlant / dewelcke deur alsulcke publicque Die- ben ende Vyanden van 't Gemeene beste / niet alleene in haere neeringe becortet worden / maer oock de schade ende 't achterheydt van 't Gemeene by den selven veroor- sacckt



Categorisation

- Manually labelled ± 3000 laws.
 - Used in the casestudy: 470.
- Hierarchical categorisation – 4(+1) layers (MPIeR).
 - Converted to SKOS format (Short Knowledge Organisation System)
 - <https://doi.org/10.5281/zenodo.3564586>

Categorisation



See <http://journal.dhbenelux.org/journal/issues/002/article-23-romein/article-23-romein.html>



Training Annif

SET-UP

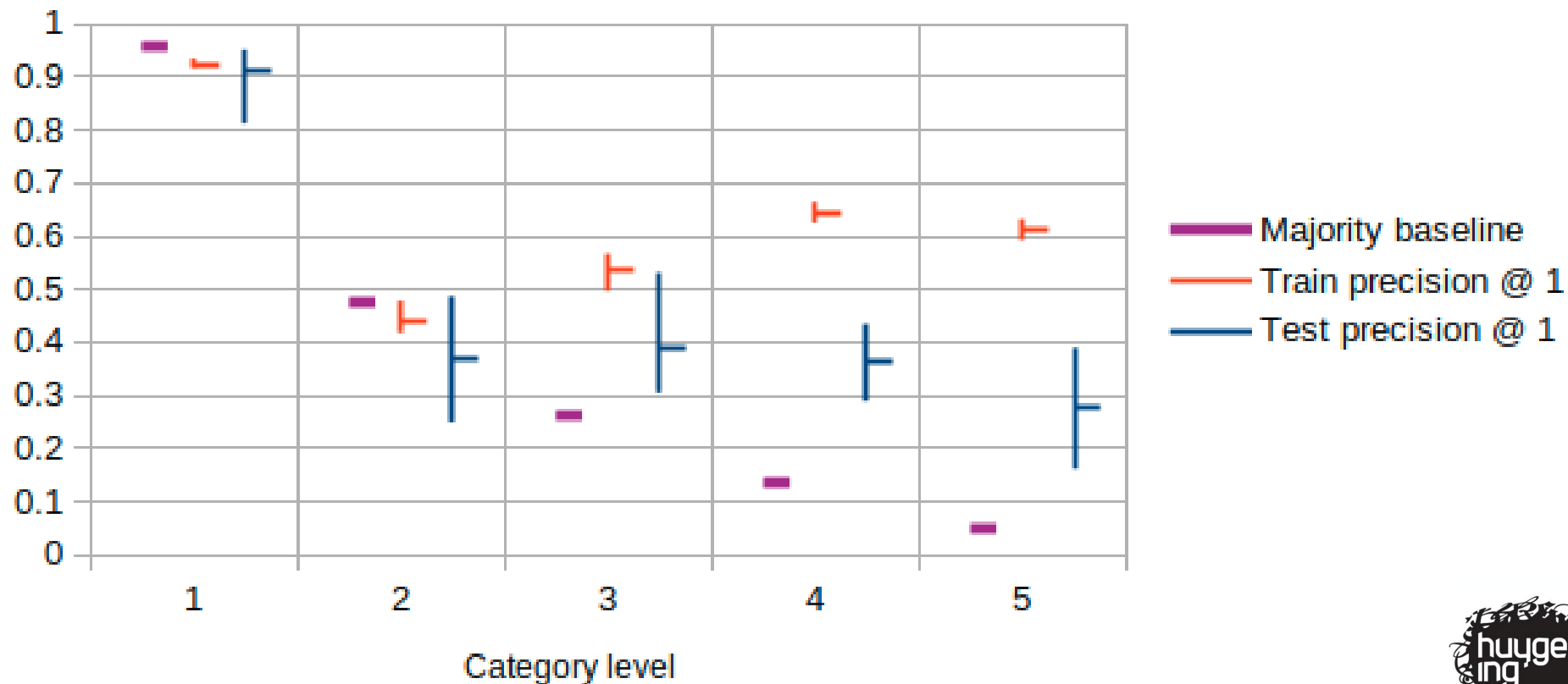
- Proof of concept: 470 labeled 'documents' (laws)
- <10 subjects each, 3.3 on average
- 69% and 28% of annotations at level 5 (deepest) and 4

- Training per category level
- 10-fold Cross Validation
- Backend: TF-IDF
- Limit = 4, threshold = 0.4

Training Annif

RESULTS

Precision @1

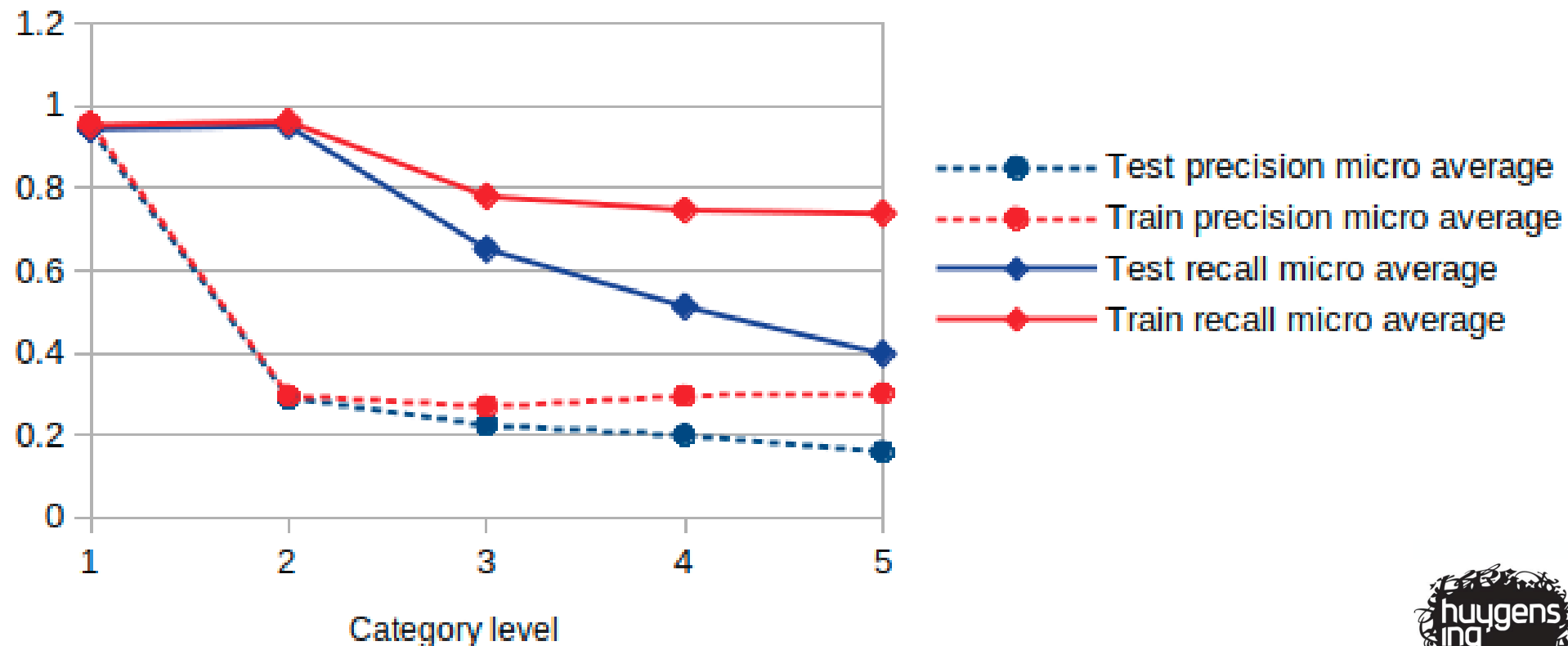


Training Annif

RESULTS

Micro averaged recall & precision

Limit=4, Threshold=0.4



Training Annif

DISCUSSION

- Limited amount of data
- Only 'dumb' TF-IDF – curious about other backends
- Impact of OCR quality?
- Stemmer: Snowball Dutch
- Impact of historical language & spelling variation?
- Making use of the hierarchy in categories

Further Reading

POINTERS

- DH Benelux journal publication:
<http://journal.dhbenelux.org/journal/issues/002/article-23-romein/article-23-romein.html>
- Info on KB lab:
Dataset: <https://lab.kb.nl/dataset/entangled-histories-ordinances-low-countries>
Blogs: <https://lab.kb.nl/about-us/blog/entangled-histories-bumps-road-and-bursts-success>
- Other work with Annif at KB:
<https://zenodo.org/record/3899723#.X2B2N1bRZJd>



Further Reading

FUTURE WORK

- Annemieke continues to work on early-modern legislation (<https://en.huygens.knaw.nl/projecten/game-of-thrones/>)
- At the KB, we're experimenting with Annif in a tool to aid the cataloguing department: 'Demosaurus'.

See <https://zenodo.org/record/3899723#.X2B6Q4bRZJc>



KB } nationale
bibliotheek

