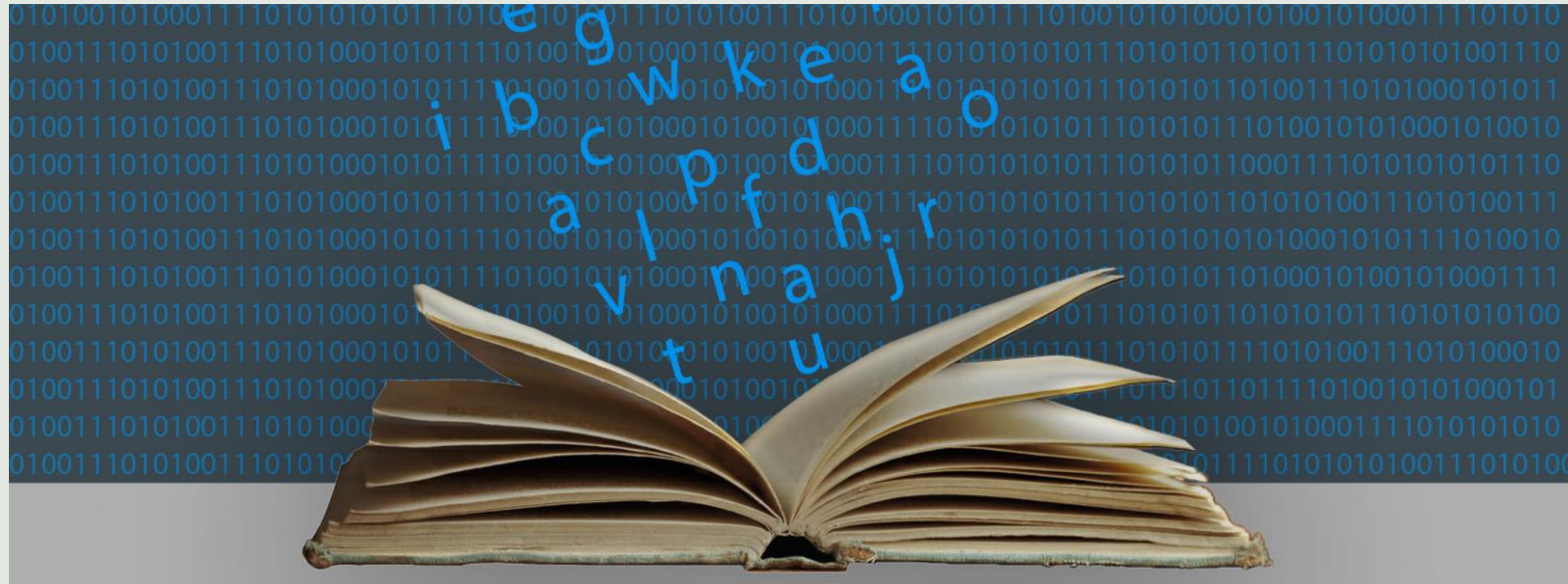


OCR quality matters in Digital Humanities research



Mirjam Cuper, 1 maart 2021

Research

- Poland
- 1672
- Oprechtse Haerlemse Courant



ITALIEN.

Geous den 3 Augusti. De Wapen-plaets van onse Troupen is geordonneert in Albenga, om de Savoyaers daer te stuten, welke, in twee Esquadrons verdeelt, vermeynden haer t'amen te conjungeren, omme Albenga te attacqueren; maer is belet geworden van den Colonel Ristori met sijne Corfi; doch also sy alle Krijghs-Instrumenten in Mombaldone hebben gebracht, ducht men, dat sy Albenga sullen willen belegeren. Dezer dagen waren eenige Troupen, soo Cavallery als Infantery, uyt Oneglia getrocken, stroopende in Diano en Cervò; doch 3 van onse Galeyen en 600 Corfi van de kant van Arasli daer by ghekomen zijnde, deden haer weder retireren.

Turino den 4 Augusti. Donderdag laest quam hier een extraordinair Courier van den Paus met een Brevet, aennemende onsen Harogh van sijne Differenten met die van Genova door Arbiters af te doen. Ondertusschen heeft men een ordre gepubliceert, waer by alle Genouesen belast werdt haere Goederen van Villa Franca en andere Zee-plaetsen wegh te halen, en daer niet weder te komen: en men wear alle mogelijke vlijt aen tot versterkinge van onse Armee, die reets 25000 Combattanten soude sterck zijn, behalven 60 Compagnien van Ordonnantie, die men mede sal oeffenen tot dat Empey.

Roma den 6 Augusti. Den Paus, beyreeft dat den Savoyen Oorlog gants Italien sal ontstecken, heeft geresolveert een extraordinair Nuntius aen alle de Princen van Italien te senden, ten eynde, dat die haer te samen souden willen verbinden tot bescherminge van de Vryheyt van Italien. Ondertusschen schijnt den Paus maer tamelijck te vreden te zijn met de Electie van Mousr. Spada, als noch te onervaren, omme sodanige geschillen, als die tusschen Savoyen en Genova, af te doen. Den Paus, met dese extraordinair hitte veele Frygten gegeten hebbende, vint sich wat onpasselijck van een sinckinge op sijn Been.

Milan den 7 Augusti. Onsen Gouverneur, voorsiene eene onverwachte overval, is Saterdagh laest gaen besichtigen op de Plaetsen op de Tesino, om aldaer die te doen fortificeren, welke tot onse verseeckertheyt mogen dienen. Van 6 Compagnien, nae Final gesonden, om Schepen te gren op de Galeyen van Turis, zijn de meeste Soldaten overgelopen tot de Genouesen, die door goet Geldt op de hande veel Volck aen haer trekken.

Lvorno den 8 Augusti. Een Flotta van Genova brengh abt dat de Savoyaers, boozstroopende gants Diano en Angnaglia, een Schip met de Genova sal was ontsaen / daer in veel Weith / son han d' ren sig d' ander kant / was gebleven; en dat de Genouese Galeyen een Frans Scheppe van Oneglia met Onse hebben genomen. Op den 5 deser arriveerden hier 't Schip la Billa nella Gitta / Capiteyn Carbonella van Alexandria / hebdenke onderwegen verstaen dat de Maltheesse Capers een Schip is met 40000 Stucken van Wechten hebben bekomen: en op heden is hier gekomen 't Schip d' Landjacht / Capiteyn Emis Engelsman / van Turis in 26 dagen / met de confirmatie van de Orde van de Franse / helgheus welcke 500 Franse Slaven wederen weg gegeven / sommege booz Selbt en andre booz niet; en dat de Galeyen van Biserta op den Hoof uyt waeren. Op heden sijn hier in 't aensicht gekomen de 13 Schepen van 't Hollandsche Convooy van Amiens.

Venezia den 12 Augusti. Tot Gomenizza is een grooten Brandt geweest. Met een Schip van Genova heeft men, dat 2 Galeyen van Biserta hadden genomen een kleyn Scheppe, wille...

dat Mr. de la Haye met 10 Franse Schepen was gearriveert aen Suratte; van waer mede genomen hadde Sr. Caron Directeur Generaal van de Franse Compagnie in Indien; en dat vervolgens was vertrocken nae Goa, van waer noch 5 Galjoenen hadde mede genomen, so dat omtrent 20 Zeylen sterck was, waer mede soude gaen attacqueren Conchin; en 't schijnt, dat men hier de Oost-Indische laecke seer ter harten neemt; weshalven men op alle wijzen de Portugyen tegens de Hollanders tracht op te maecten; zijnde onlanghs, soo men segt, omme Gelat-middelen te vinden tot dien eynde, tot Lissabon een samenkomst geweest. Hier is men belig met noch meerder Krijgs-macht op de Been te brengen; weseude de Tydingen uyt Spangie seer twijfelsachtig, of wy Oorlog of Vrede met die Kroon sullen hebben, alsoo self den Spaenssen Raedt, soo men hier segt, niet eens was. Daer is hier 100000 Kroonen van den Onsfinger Generaal van Parijs ge-cyft, tot een Fonds voor de Leger-plaetsen van de Generaliteyt. Mr de Gaumont is na Turin gesonden, omme Duc de Savoyen te bewegen, dat sich met die van Genova soude willen accomoderen.

POOLEN, PRUYSEN, &c.

Warschou den 9 Augusti. De Coninglijke Reyse gaet aen slaende Donderdag voer: sijn Majest: sal de Koningin Maria in de Goedeten van den Bisschop van Polen houden: En alsoo de Tydingen uyt de Ukraine vast alle dagen gevaerlijcker aankomen, gaen nieuwe Unversalen af, ten eynde, den Adel te eerder te doen opsluten. De Turcken souden niet meynen in de Ukraine hiet op te houden; maer voort in Poolen in te vallen, omme dat te verwoesten.

Danzigh den 13 Augusti. 't Laet sich in Poolen wonderlijck aen sien. 't Gevaer schijnt groot, en nochtans komt den Adel traegh op. In tegendeel maecten die van de Franse Factie zig geheel sterck, om haere Partye noch te doen prevaleeren.

Duytslandt en d'aengrensende Rijcken.

Weenen den 11 Augusti. De Evangelisse Predicanten zijn van Presborg vervoert, en sullen noyt meer in de Erflanden mogen komen. Den Spaenssen Courier is weder te rugge ge-expedieert, alsoo de Keyserlijke Troupen nu alle volkomen in marsch zijn.

Weenen den 14 Augusti. Gisteren heeft den Heer Tr. Generaal Montseruill sijn afscheyt Audientie by sijn Keyserlijke Majest: seep gehad; en vertreckt daer op heden met sijne gantsche Camp nae Oger / zijnde den Prince van Lottharingen vers derwaerts vertrocken: men segt de Keyserlijke Troupen geselt tot Wittenbergen sullen marcheren; omme met de Keyserliche Troupen in Sultistat-landt te haeten; daer gaen der Keyserliche en Duytse Caballiers inde booz Bolontaire. Den Generalen Compagnie Giersteins is by de Keyser aen sijn in 't aensicht gekomen. In 't Geseelschap booz sijn de Duitzen met 10 Mill ver booz den Dender geseleert.

Minden den 15 Augusti. Het is geheel seecker, dat sijn Cheurvorlijcke Doord: den 22 deser van Berlijn sal opbreecken; de Velt-tuyg-meester Goits is reets voor af vertrocken met 6 a 7 Regimenten, om 2 Mijlen van hier sich te conjungeren met de Troupen van den Keyser. Onsen Cheurvorst heeft al sijn Volck uyt Pruyssen op ontboden; en wert gelegd sijn Armee effectief 32000 Man sal sterck zijn, behalven die van den Keyser, Denemarcken, Saxon, Hessen en andere. Den Francken Ambassadeur doet groote eforten, om onsen Cheurvorst om te setten, presenteerende hooggedachten Cheurvorst te restituieren alle sijn Vestingen, en daer-en-boven wel 10 Tonnen Gouts; maer by vint geen

Research



POOLEN, PRUYSSSEN, &c.

â- iÃ- AVarfÃ©J:Â«iuÂ»>,p-f-',,aij, Mr>'. 'â€¢, f ua-jctt
teJ>!r&%FjÃ- b â- 'Â» t^ptgaa aÃ©-i; â- yruy-ÃªÂ« Aijai
a/bi*^; ;'vi* je-Â« .]Ã«%Â«})'Â«fcWftim s,ijcii"(dat"',Ha* ata/ifia
Â«tadifc PlaK.fcn..lijcf onder jijÂ» Mijcftcv-ji bip
;ftkMÃ»IÂ«Jp'Â»Â»i.^ ; toegrvaJlw!..f Da*!'iija defer dagi-ri
bietbÃ³ek vctfrh-y-ie 'Taf'ie'Â» [â- â€¢ dÃ©Â» ColieUr.a
s'tDgcfc'iitjeÃÿ j", aiWevtiÃ©im.'flÃ«worÃ©ftic
!WÃ«ptfH*ei:krp"i'a,t,^ u'tenj'^Â»K*Site^^),epÃ^r;oÃ³d<Xi^ â-
'â- : <â- ,; . "Ã-i'iÃ- 'boV-l 'Â«.'â- ."â- a*Iri" '<Cbr-if frgft

grnfsan â€¢jj* . ;rf*tt
fflÃ- i.b-'rf/it.aÃ- ifie
'teÃ©i:sSÃ- 4te (jic
oÃ»-Â»Â»

â- 'Ã- Â»n'iW|iteiiÃ

POOLEN, PRUYSSSEN, &c.

' Wurfihotiw den 6, February.'VyÃ- dÃ" UkÃ- faiiiÃ© kiÃ³men va-"
rÃœblÃ"Tydiagea, * dief iwÃ©fdt'gelchreweitj"dat'dÃ«nDorÃ¶fehsko
ganti blbÃ¶tÃ©nalÃ£ iÃ«er k:ha<rs';..Ã«ii de: Padden raÃ«eit
weghViwaerorn de CÃ¶f hÃ«tri teVerlatca, "Ã©n haer nae" deÃÿ
Hanenko'tÃ«begeyca ,die defwÃ«genhydeiooÃ¶¶¶ Man cflÃ«ftivÃ«'
IÃ³udÃ© fterck"Ã¼jn C, waer vaa "inÃ«adÃ© confirmatie
vÃ«nvacht.*i-t'; ;'..;â-

POOLEN, PRUYSSSEN, &c.

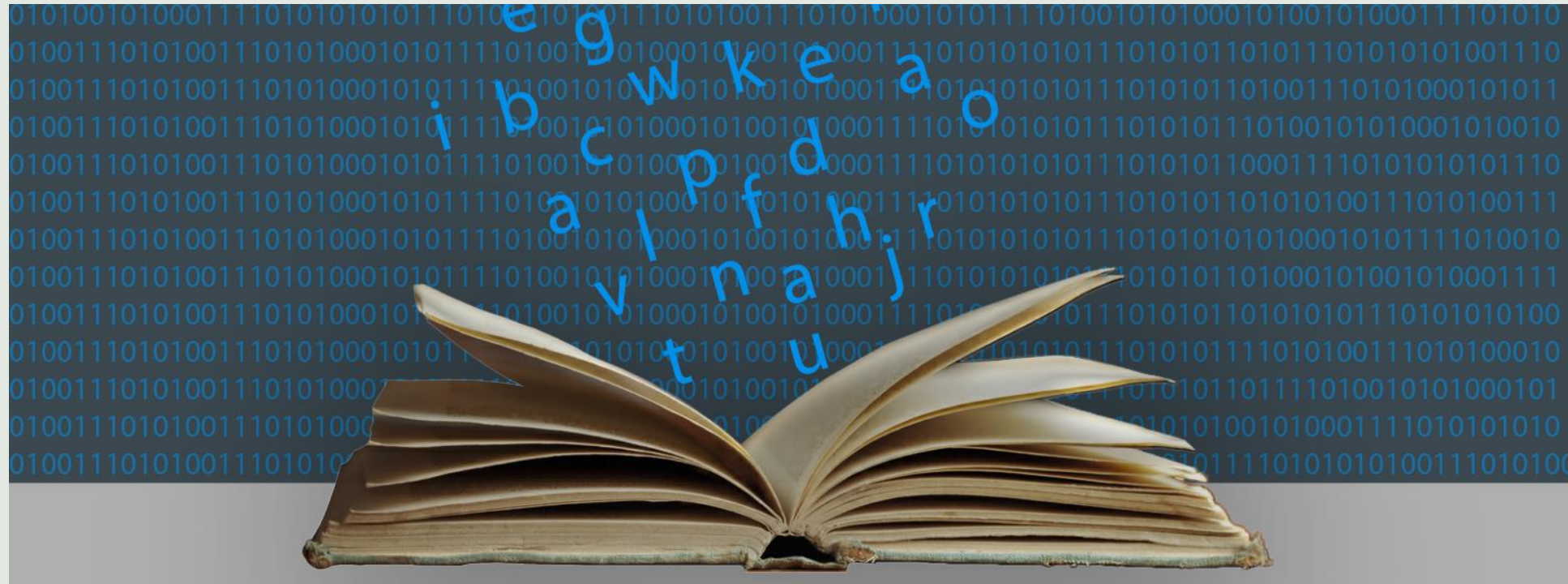
â- iÃ- AVarfÃ©J:Â«iuÂ»>,p-f-',,aij, Mr>'. 'â€¢, f ua-jctt
teJ>!r&%FjÃ- b â- 'Â» t^ptgaa aÃ©-i; â- yruy-ÃªÂ« Aijai
a/bi*^; ;'vi* je-Â« .]Ã«%Â«})'Â«fcWftim s,ijcii"(dat"',Ha* ata/ifia
Â«tadifc PlaK.fcn..lijcf onder jijÂ» Mijcftcv-ji bip
;ftkMÃ»IÂ«Jp'Â»Â»i.^ ; toegrvaJlw!..f Da*!'iija defer dagi-ri
bietbÃ³ek vctfrh-y-ie 'Taf'ie'Â» [â- â€¢ dÃ©Â» ColieUr.a
s'tDgcfc'iitjeÃÿ j", aiWevtiÃ©im.'flÃ«worÃ©ftic
!WÃ«ptfH*ei:krp"i'a,t,^ u'tenj'^Â»K*Site^^),epÃ^r;oÃ³d<Xi^ â-
'â- : <â- ,; . "Ã-i'iÃ- 'boV-l 'Â«.'â- ."â- a*Iri" '<Cbr-if frgft
grnfsan â€¢jj* . ;rf*tt]jÃ£'fcuteÃ«itf'te &^^ ; ftefsÃ»' teiX-S fÃ¶ fflÃ- i.b-
'rf/it.aÃ- ifief fyoao Jfrari'fy f-Ã©fC;WtPa*t *epa*ji 'teÃ©i:sSÃ- 4te
(jidf sa^teteintn;]nuÃ«bimiJ.Â«'Ã¼cÃ³jaant,-; oÃ»-Â»Â»

â- 'Ã- Â»n'iW|iteiiÃ- *swj.SÃ«'iSiÃ«B'tTf

POOLEN, PRUYSSSEN, &c.

. S Warfchonm 'itol Den __t Â«cat dagen gsduu-t heeft, Ã¼ eynA-
ngg^aoldgtheaÃ³ymds â- "V-â€¢â- ,<â- ,'.", ',â- ., â- ':â- â- -,â-
TT-','; ;"r-r â- â- â- "-â- ': : â€¢â- â- ' -\ : - .â- .â€¢.: -gTdoopons;
iaÃ- g?HjckÃ« . heefti_oOj:fc dn, t.fatT)enkon.ft. VfanSaiq; domirfch
toeOjjataw i^h^ewÃ«nfchtÃ©aeyn_bciomen..!Da TÃ¼rekft-liier
btea Sm Â«ccl te weg? gciarâ€žfit<.dat."dÃ©n.. Adel;
% Ã³cr.onaÃ- nrghwasiâ€ž nitmar_illÃÿ.'andflr-Ã- :
rzydeo geÃ« ft* v*n-dÃ»VSDirland:BidefemiÃ¶ tÃ©
ge*. li_j^fcÃ£lfcmm4ij_c&'IÃ- ecftTviÃ- a^ Voeten40000
l verrftsrektwerden /': Deri Heer Roi zfesky.P_iÃ£lijckcn
' , en voordÃ©GoningitineÃ©ea RoiÃª mbd9llreaght,â€¢
Ã© tÃ³rcyts tot GratÃ¶aw'.aeAgebri:' zij m â- ~
fcliein wil riiÃ«riyoÃ¶tgÃ©ldÃ³rbaÃ«rberichten., ditden
miriEoÃ«nftcett6 l.cnKtiMiPoolenfilk'OJnea, en tot van
Zywiec falVÃ©rtrecke," t wÃ«kk
en }â€¢ doch hier van vernsichtinÃ«n
>. ly ;...-''â- ; "â€¢ - -â- ; 7,7' 1 .â- -

OCR quality matters in Digital Humanities research



Mirjam Cuper, 1 maart 2021

Introduction

- Mirjam Cuper
- Data scientist
- KB - National Library of the Netherlands
- Digital heritage



DIGITAL HUMANITIES

Humanities

- Human society
- Human culture

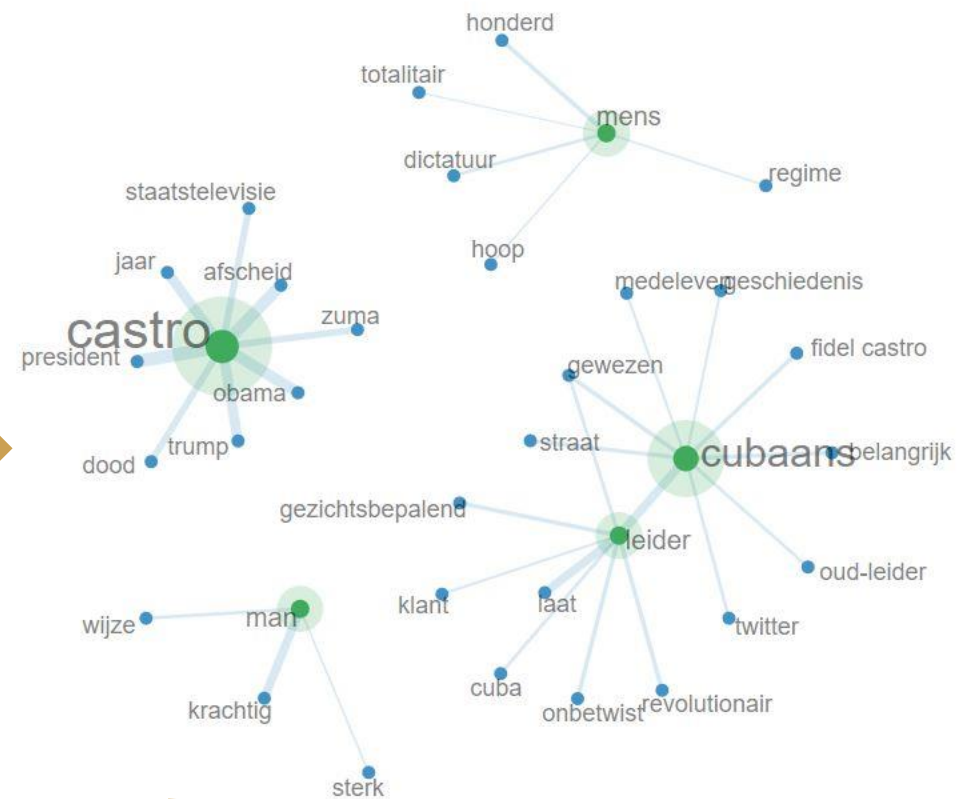
- Heritage collections

Optical Character Recognition (OCR)



Cholera in Peru uitloper van 'pandemie'

ATLANTA (Reuter) — De cholera-epidemie in Peru wordt veroorzaakt door de El Tor-stam van de ziekte, die voor het eerst in 1961 werd geïdentificeerd in Indonesië. Dit heeft het Amerikaanse overheidscentrum voor ziektebeheersing CDC in Atlanta gisteren bekendgemaakt. De El Tor-bacterie heeft zich in de afgelopen dertig jaar bijna over heel de wereld verspreid in wat het CDC een 'pandemie' noemt. Het instituut merkt op dat Zuid-Amerika tot afgelopen maand gespaard was gebleven.



Genre classifieer

Churchill, die jarenlang de eerste plaats innam onder de meest bewonderde mannen, doch deze het vorige jaar moest afstaan aan dr Drees, heeft kans gezien de eerste plaats weer te halen met een stijging in populariteit van vijf procent, terwijl dr Drees met een daling van vier procent weer op de tweede plaats kwam te staan. Dit is een van de resultaten van het onderzoek van het Nederlands Instituut voor de Publieke Opinie, dat onlangs alleen aan mannen over het gehele land verspreid en uit alle lagen van de bevolking de vraag stelde: Welke van alle nu levende mannen, leden van de koninklijke familie niet meegerekend, bewondert U het meest? Churchill kreeg vijftien procent van de "stemmen", dr Drees 13, Eisenhower 6, Jan van Zutfen 6, Albert Schweitzer 3, Paus Pius XII 2, Adenauer 2, oud-rinifter Liefstinck 2, Einstein 1 en Abe Lenstra tenslotte ook 1 procent. Jan van Zutfen, Adenauer en Einstein komen dit jaar voor het eerst op de lijst van de meest bewonderde mannen voor. (ANP)

Nieuwsbericht: 87.87%

Column: 1.55%

Achtergrond: 3.55%

Reportage: 0.77%

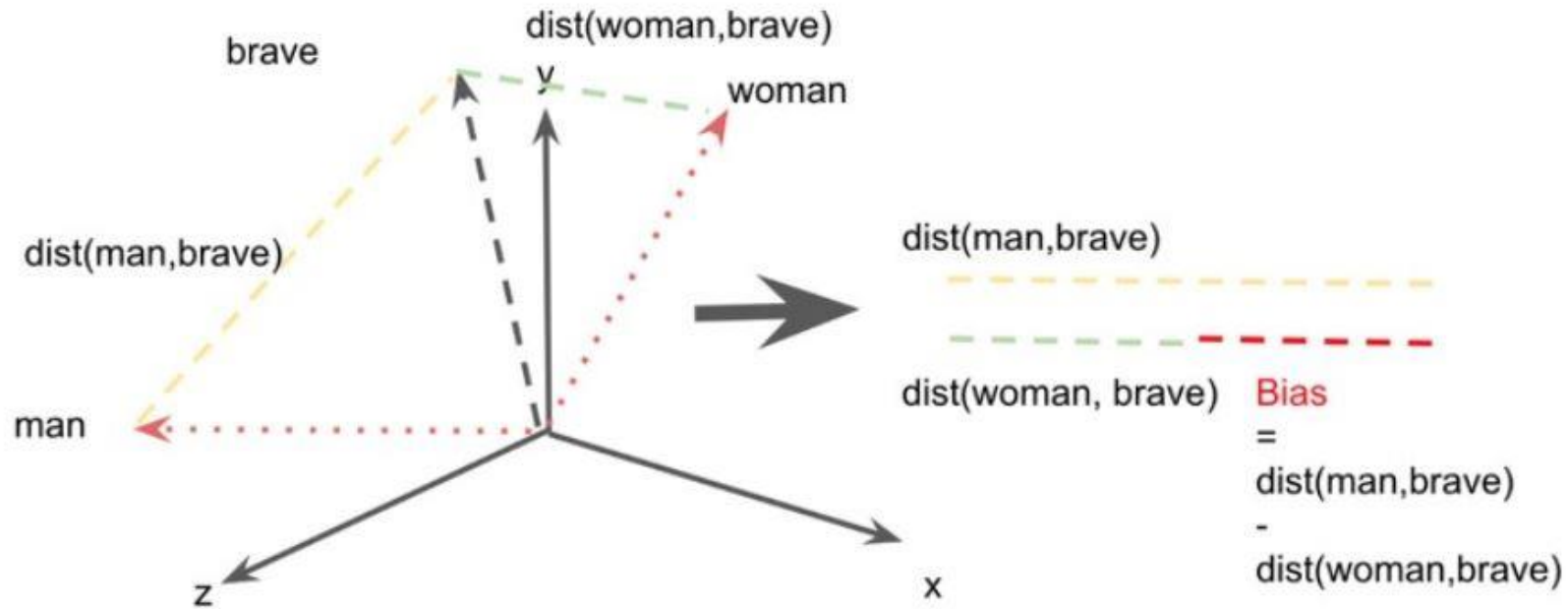
Verslag: 3.28%

Recensie: 0.81%

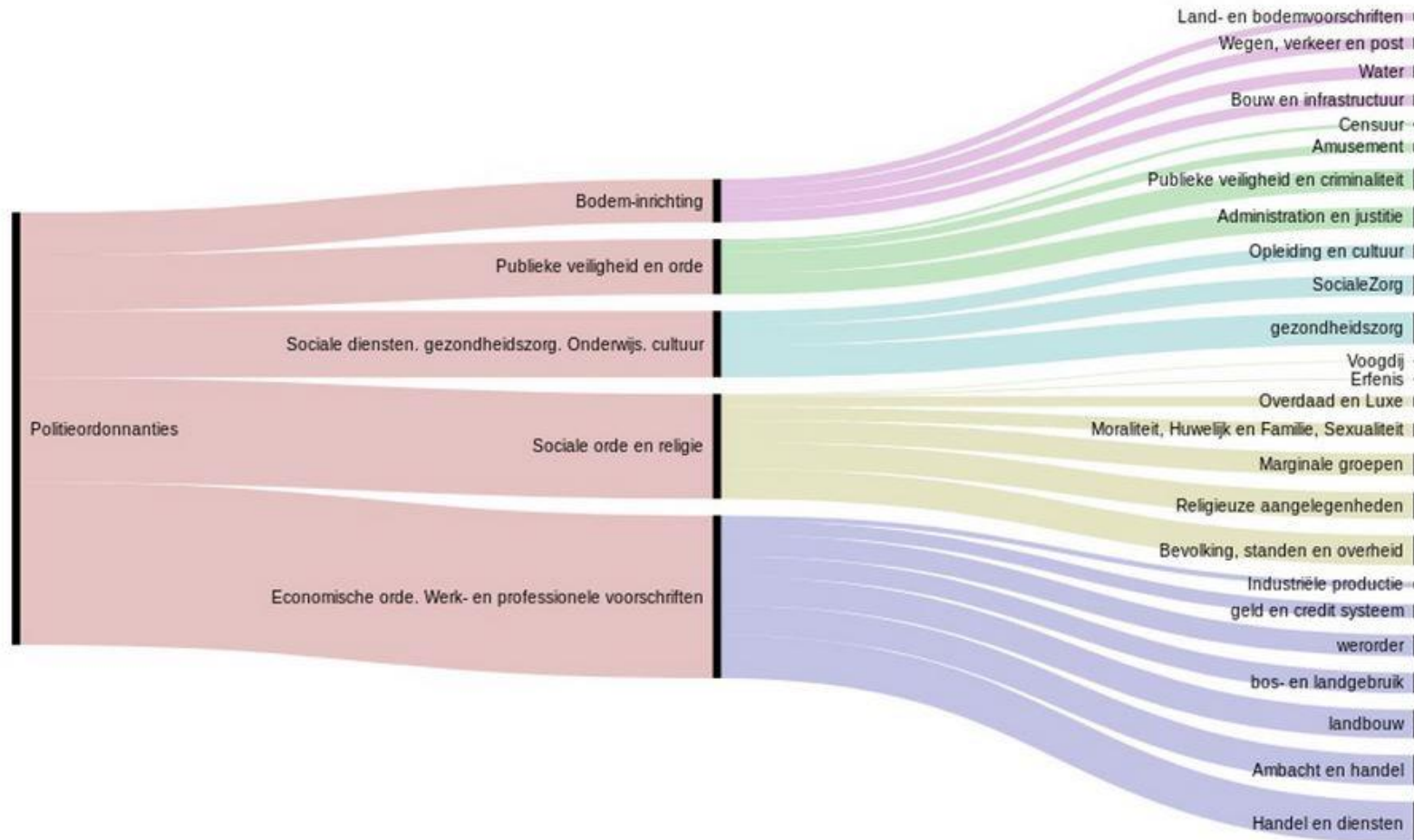
Opiniestuk: 1.91%

Interview: 0.21%

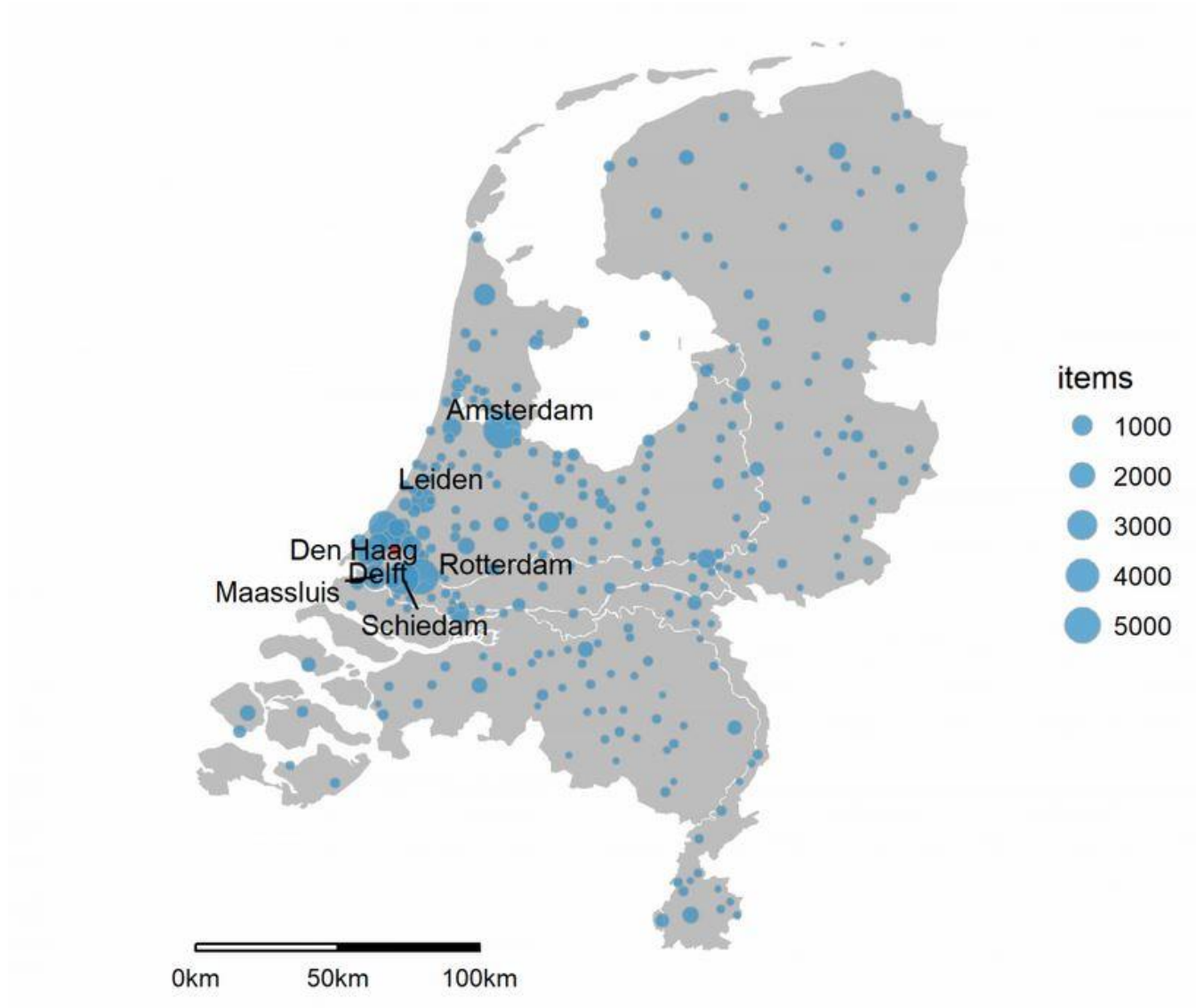
Gender bias in historical newspapers



Classify books of ordinances



Cities in historical newspapers



OCR PROBLEMS

INFORMATION RETRIEVAL

Information retrieval

Kranten ▼ vrouw Zoeken Uitgebreid zoeken ▼

[Terug naar kranten](#) **4.676.780 krantenartikelen gevonden**

Zoekterm vrouw ×

Sorteer op relevantie ▼ i Weergave 📊 ☰ ☰

Periode


- 17e eeuw (38)
- 18e eeuw (7571)
- 19e eeuw (463098)
- 20e eeuw (4178186)
- 21e eeuw (27887)

[Kies periode...](#)

Soort bericht

- Advertentie (792330)
- Artikel (3614872)
- Familiebericht (235218)
- Illustratie met ondertekening (34360)


Verspreidingsgebied

 **CANDID TAPE Een rotzooitje** ☆

vrouw dan, waar we ons geld voor betalen? Eerste **vrouw**: Nou! Tweede **vrouw**: Waardeloze rommel Eerste **vrouw**: Waardeloze rommel.

Krantentitel NRC Handelsblad
Datum 17-04-1979

[Meer details](#) ▼

 **Openbare Leeszaal en Bibliotheek.** ☆

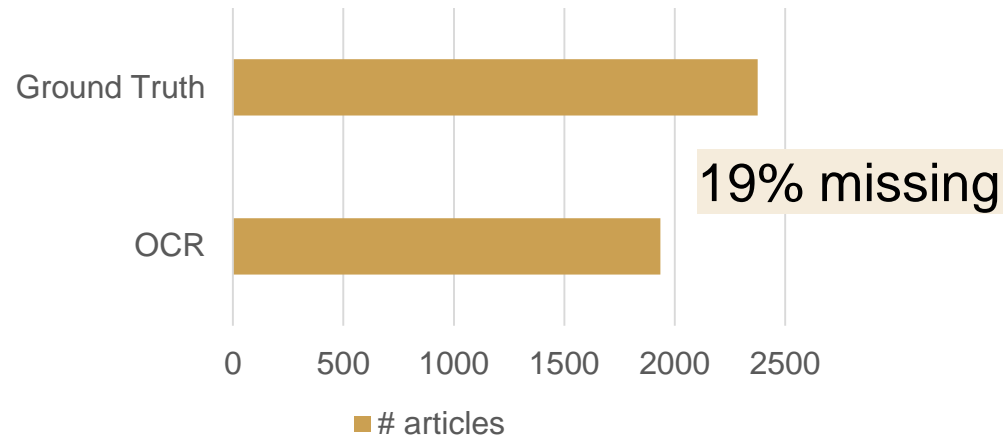
vrouw0; Een meisje-studentje" d. A. Salomons. 1907. Worp, J. A» Een onwaardelijke **vrouw**, brieven en verzen van en aan Maria Tesselschade. (Andere bladen worden vriendelijk ver

Krantentitel De Zaanlander
Datum 04-05-1927

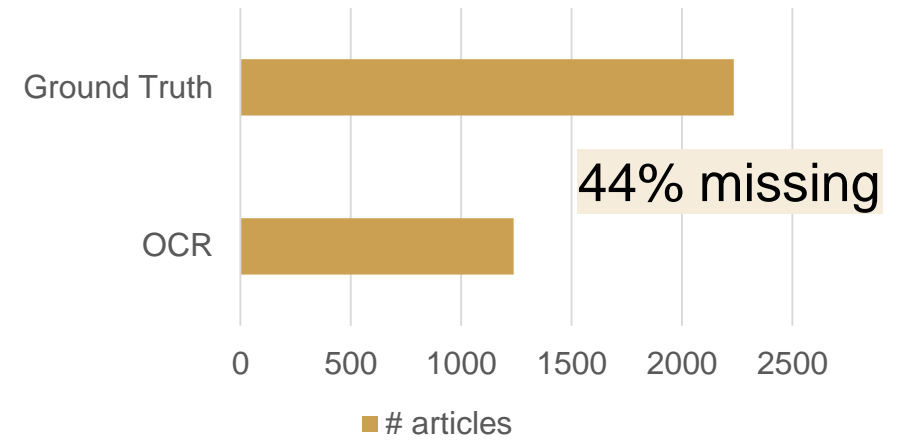
[Meer details](#) ▼

Information retrieval

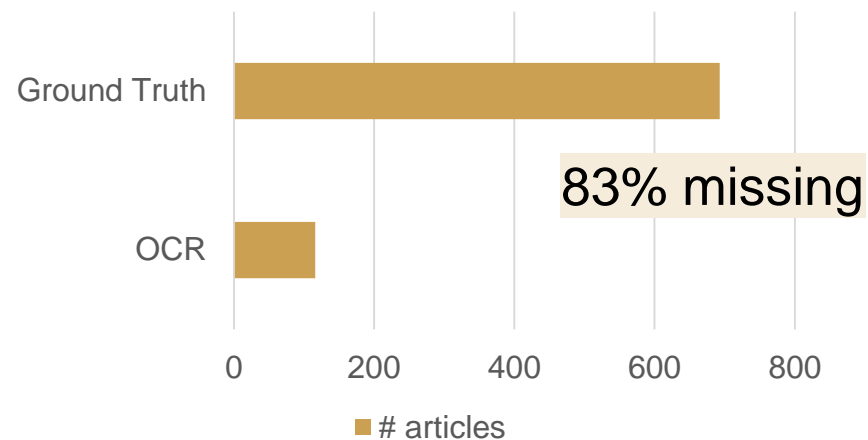
Search query 'vrouw'



Search query 'Duytsland'



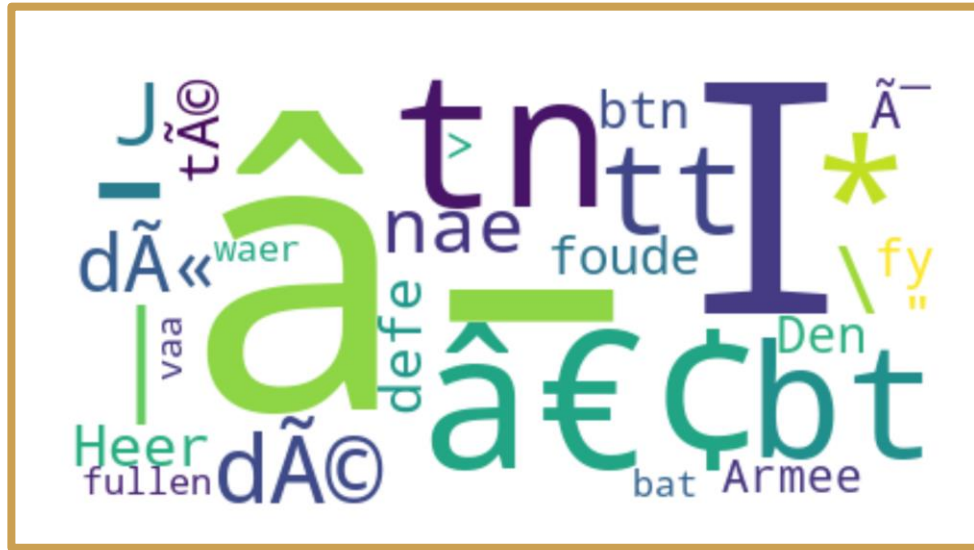
Search query 'Coninglijcke'



ANALYSES

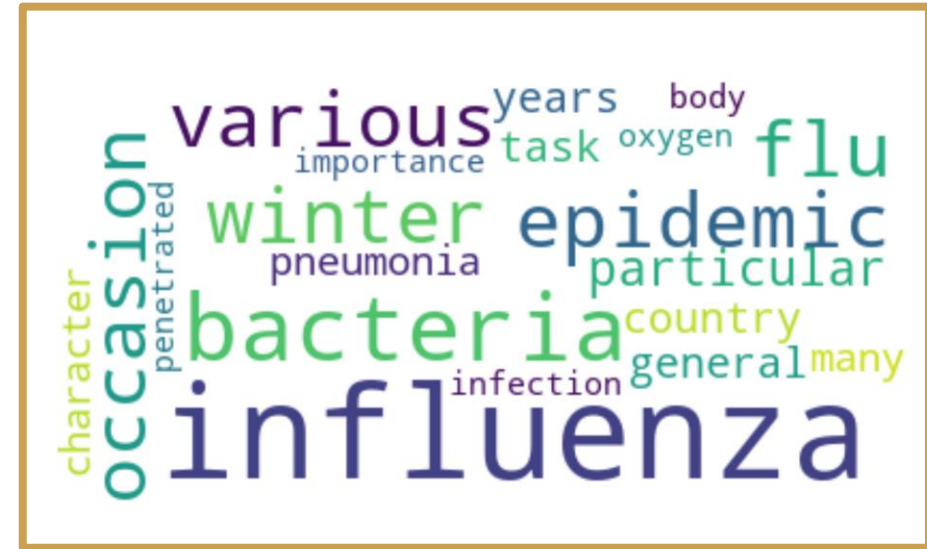
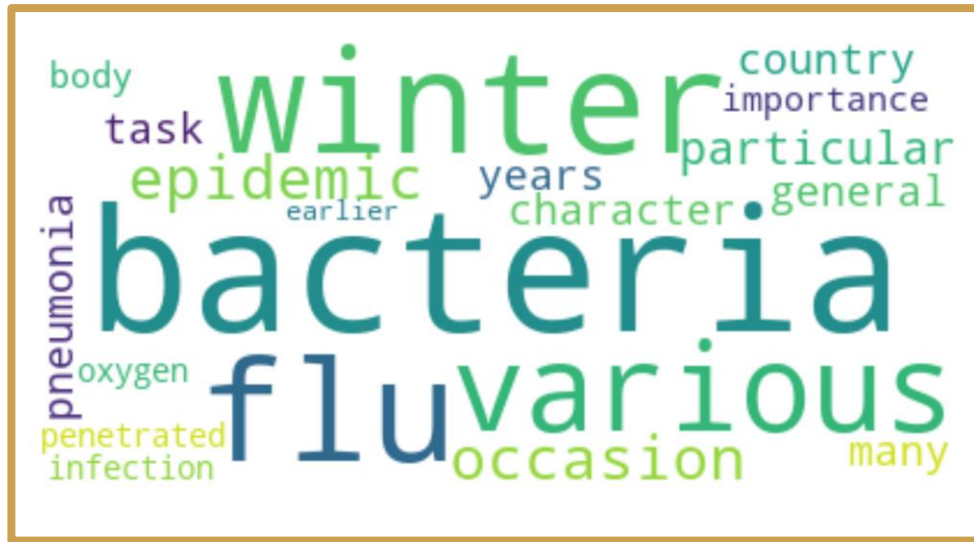
Analyses - wordclusters

Topic: reports about Poland in 1672



Analyses - wordclusters

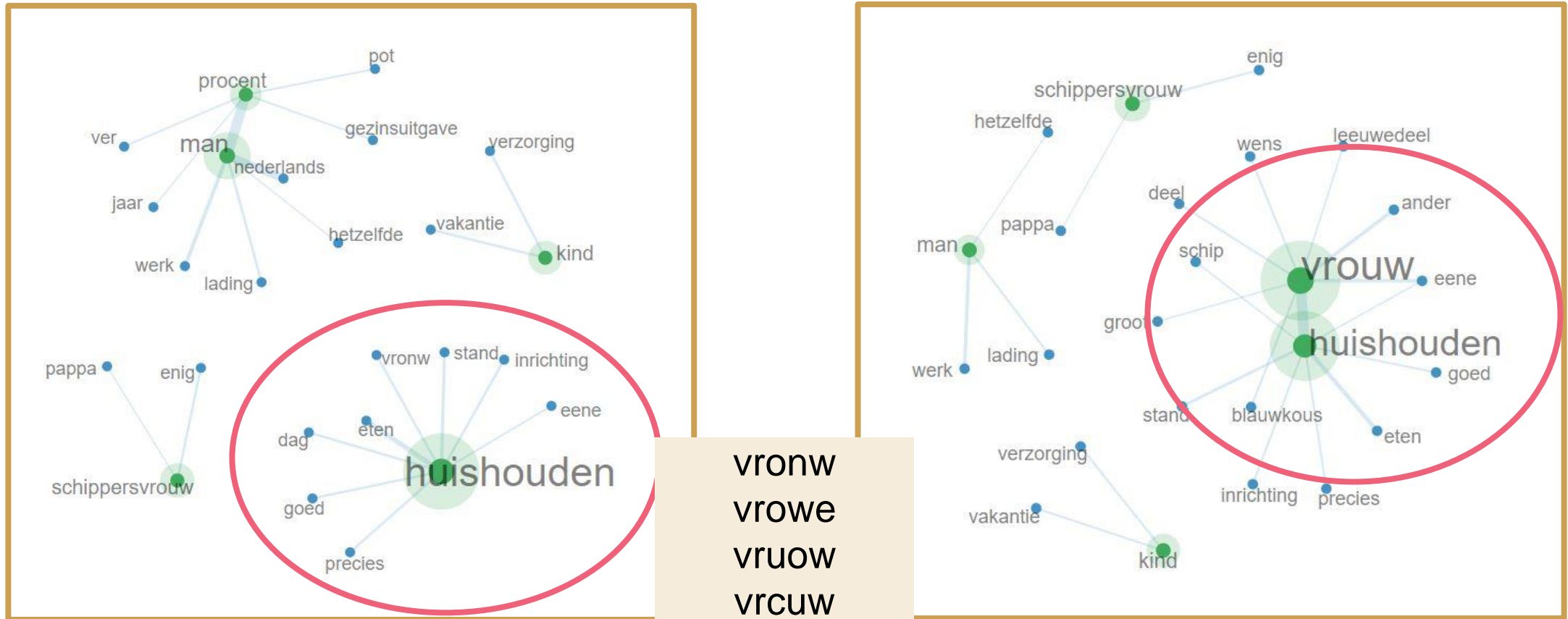
Topic: relation between 'disease' and 'influenza' in newsarticles



influenza
influemza
infleunza
imfluenza
influanza

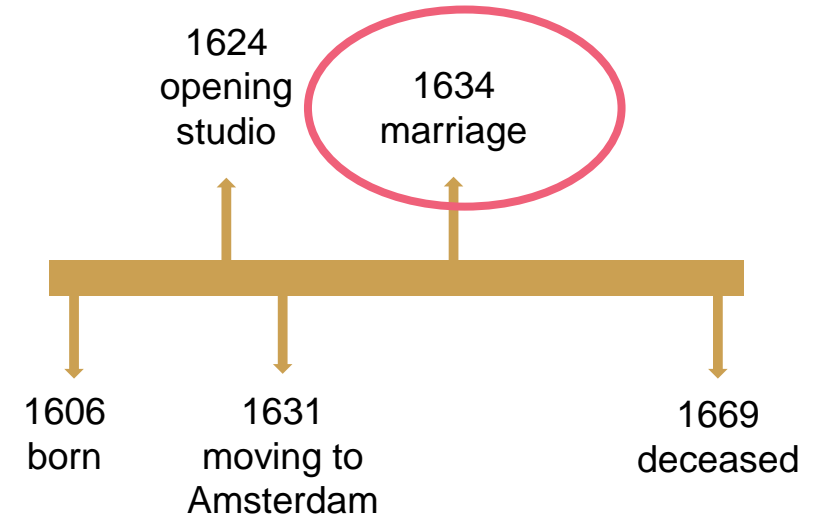
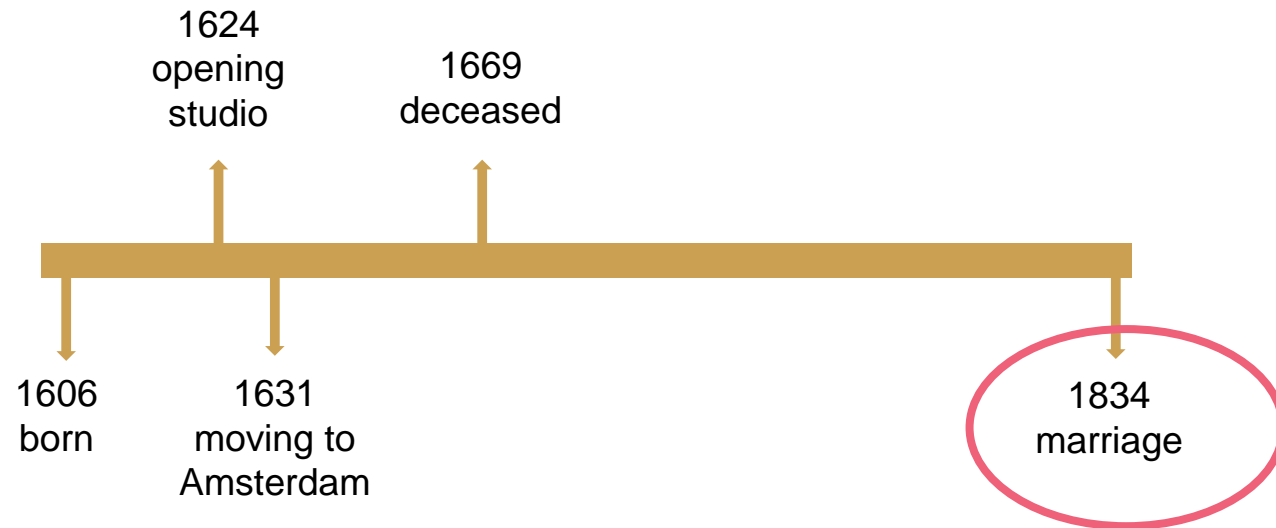
Analyses – keyword association

Topic: relation between 'household' and 'woman' in newsarticles



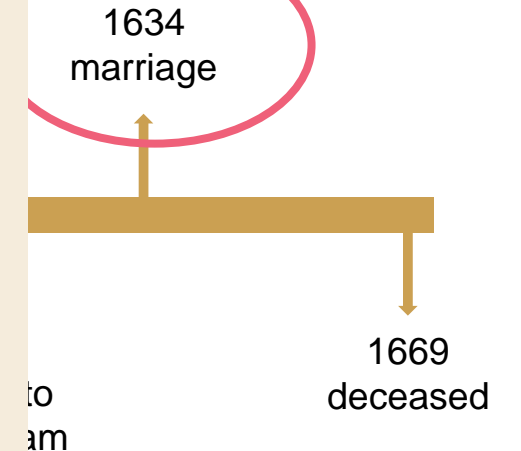
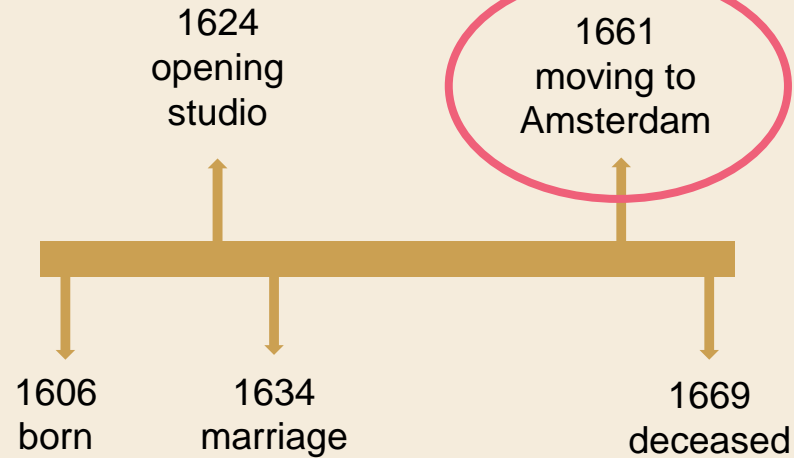
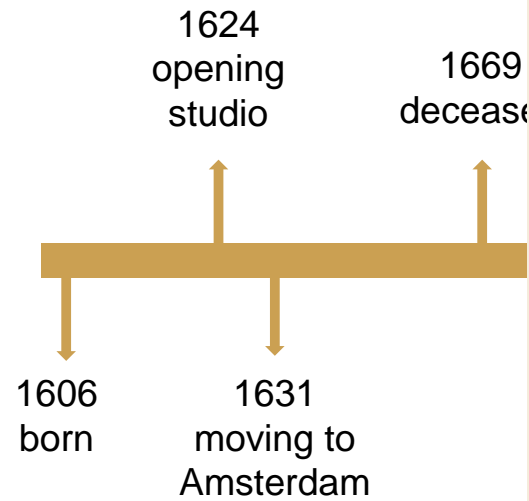
Analyses – timelines

Timeline Rembrandt van Rijn – Dutch painter



Analyses – timelines

Timeline Rembrandt





Analyses – Named Entity Recognition

Correct

Kroonprins Willem-Alexander PERSON bracht dinsdag DATE een bezoek aan Zeeland GPE . Hij deed onder meer het Zeeuws Museum ORG in Middelburg GPE aan. In de tapijtenzaal kreeg de prins uitleg over de Zeeuwse LOC wandtapijten.

Errors

Kroonprins Willem Alexander PERSON bracht dinsdag DATE een bezoek aan Zeeland GPE . Hij deed onder meer het Zeeuws Musuem ORG in Middelbvrg GPE aan. In de tapijtenzaal kreeg de prins uitleg over de Zeeuvvse GPE wandtapijten.

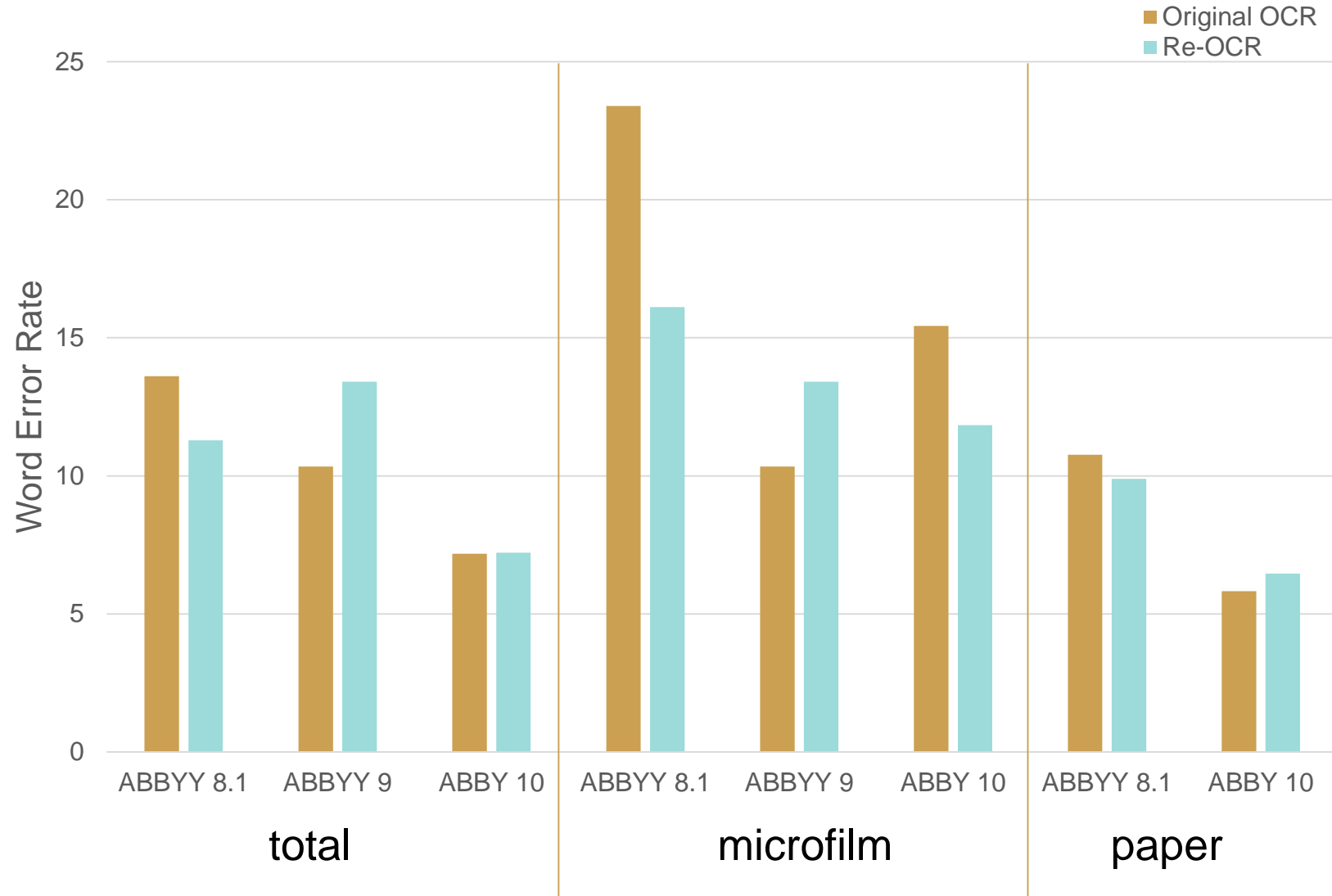
WHAT TO DO?

Optimize digitisation process

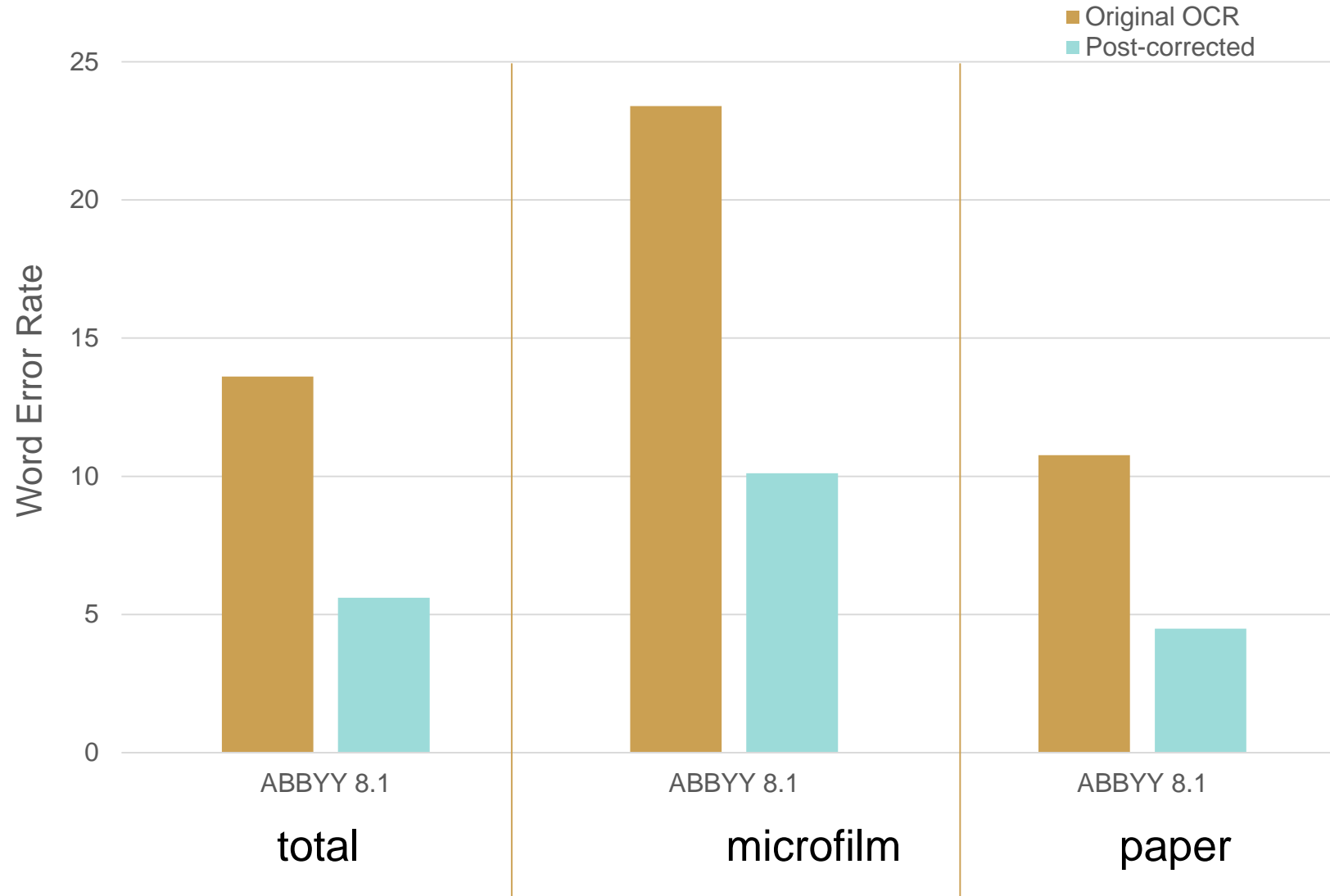
- Pre-select material
- Software based on document characteristics
- Most recent software



Re-ocr old material



Post-correction



Transparency

- OCR software
- Post-corrected or not
- Data source (microfilm, paper)

- Estimated quality



Data accessibility

- Other interfaces
- Expand queries

The screenshot shows a search interface with a dark blue background. At the top left, the label 'Zoekterm' is followed by a search input field containing the text 'influenza'. Below the search field, the text 'common OCR errors' is displayed next to an information icon (a lowercase 'i' inside a circle). A white dropdown menu is open, listing five suggestions, each preceded by a plus sign in a teal square:


- + influenza
- + influemza
- + infleunza
- + imfluenza
- + influanza



Improve OCR quality,
improve Digital Humanities
research!

Questions?





KB } national library
of the netherlands