# Recognizing Extracted Entities for the Historical Database Suriname Curação

Erik Tjong Kim Sang Thursday 30 November 2023





### **REE-HDSC** project



Radboud University started the HDSC project in 2017

That project aims at collecting and make available historical data on inhabitants of Suriname and the Caribbean

The task of the eScience Center is to examine how the required text-todata conversion can be automated

See: https://www.ru.nl/hdsc/



### **HDSC People**

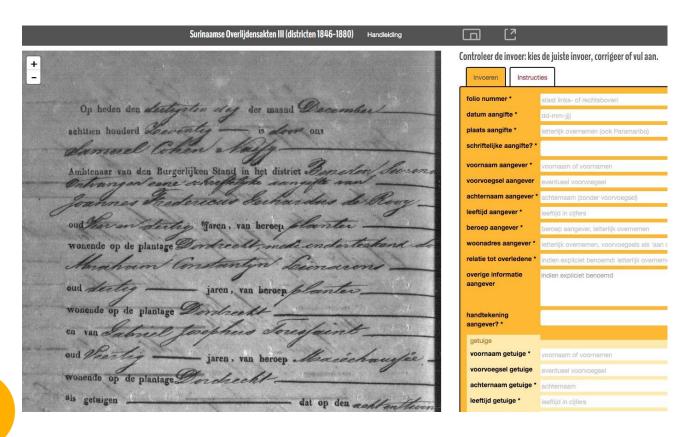
- Matthias Rosenbaum-Feldbrügge, Radboud University Nijmegen, LA
- Björn Quanjer, Radboud University Nijmegen, Researcher
- Thunnis van Oort, Radboud University Nijmegen, Researcher
- Lisa Hoek, Radboud University Nijmegen, Data scientist
- Coen van Galen, Radboud University Nijmegen, Project Manager
- Erik Tjong Kim Sang, eScience Center, Research Software Engineer

See: https://research-software-directory.org/projects/ree-hdsc

Heden den Negen en twintegsten Sanuary des Saars No En dunend acht honderd En en Dertig-Verbeteringen. Compareerden voor mij Amblenwar van den Burgerlijken stand op dit eiland. ale hanter De personen van Isaac fokannes Rammelman Elevier Aanteekeningen. Junior en Preter Hetz. van competenten ouderdom en alhier woonachtig. Dewelke aangific hebben gedaan dat of dist ciland op den tweeden fanu Dene inschuyerns is any des faars den duisend acht honderd een en dertig geschied ingevolge Bublicatie van Dine to half vefure des morgens overleden is: Pacob Bonaventura tur adm in Rade to Hammelman Elsevier 5- 13 January 1831. Oud Les maanden en twee dagen -Geboren te Curacao\_ op den dertigsten Suny des facio Een duirend acht honderd en dertig. zijnde Blank laatstelijk gewoond te Curação\_ Ongehund\_ En is deze behoorlijk geteekend na voorlezing op Curação dato utsupra. In stede van den



### **Current practice**



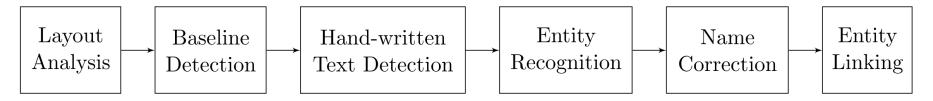
Data from scans of civil registry documents are manually entered in forms by hundreds of volunteers

https://hetvolk.org





### Proposed automated pipeline



- 1. Layout analysis: determine page layout: how many text columns?
- 2. Baseline detection: find the lines of text
- 3. Hand-written Text Detection: identify the text of the document
- 4. Entity Recognition: find the names and dates in the text
- 5. Name correction: integrate name corrections mentioned in text
- 6. Entity Linking: link names of identical persons

For step 1-3 we use the program https://readcoop.eu/transkribus



## **Evaluation of pipeline**

The hand-written text module reports a character error rate (CER) of 6.0% (94% correct), which sounds good

However, when we tried to identify names, the performance was mixed: 90% correct for dates but only 33% for names

The character error rate is always lower than the word error rate and this is again lower than the error rate for names



# Improving hand-written text recognition (HTR)

We have three ideas for improving the quality of person names by hand-written text recognition:

- 1. Additional training with only hand-written person names
- 2. Post-processing of the output of the HTR module
- 3. Automatic quality assessment of HTR output

For training we have tens of thousands of already identified names



### Name accuracy results of HTR model retraining

	1830-1839	1840-1849	1920-1929
Before post-processing	68%	69%	41%
After post-processing	68%	66%	43%
Quality estimation	93% (65%)	88% (67%)	73% (44%)

Postprocessing (replacing unknown words with known words) did not always increase performance for the best retraining configuration

High performances can be obtained for part of the data



# **HTR examples**

	Humina Martinas
The second of th	rithelmina o illis.
	Elias Judah Lon_
Enferio Vil	Slav Felisla

Name found	Evaluation
Hermina Martina	Correct
Henrietta Wilhelmina Fillis	Near miss
Maria Elias Junda Son	Hallucination
Lucina Hilda	Complete miss





## **Concluding remarks**

- Hand-written text recognition can be used for processing old documents
- This technique does not perform perfectly, in particular not for names
- Fortunately new recognition models can be trained for our data
- We trained a new model for Dutch civil records: 68% name accuracy
- Post-processing did not improve the scores
- Accuracies close to 90% could be reached for part of the documents (66%)
- Some sets of documents proved to be harder for processing than others



netherlands

Science center