# Content development efforts in the NLF's historic newspaper and journal collection: re-OCR, Named Entity Recognition and layout analysis

Kimmo Kettunen

The National Library of Finland

Mikkeli unit

THE NATIONAL LIBRARY OF FINLAND

# Background

- The National Library of Finland's Mikkeli unit has digitized historical newspapers, journals and ephemera published in Finland since 1998.

- The present collection consists of about 14 million pages mainly in Finnish and Swedish.

- Out of these about 7.45 million pages are freely available on the web site digi.kansalliskirjasto.fi (user interface in English & Swedish available)

- The time period of the open collection is from 1771 to 1929. Years 1920–1929 were opened in January 2018.

# NLF's Mikkeli Unit



- A digitization factory about 220 km north-east of Helsinki

- Ca. 40 employees

- Digitizes about 1-2 M pages of newspapers and journals annually

- Also other production

# NLF's digital humanities efforts with the newspaper and journal collection

- Besides producing and publishing the digitized data all the time NLF has been involved in research and improvement of the digitized material during the last years. We ended in July 2017 a two year European Regional Development Fund project and started another two year ERDF project in August 2017 → **Digitalia (**together with South-Eastern Finland University of Applied Sciences**)**

- NLF is also involved in research consortium COMHIS that is funded by the Academy of Finland (2016–2019) and utilizes the newspaper and journal data in its research of historical changes of publicity in Finland.

- EU Horizon Project NewsEye started in May 2018 - NLF is one of the partners and provides data for the project

# Data improvement and new ways to use the data

NLF has so far performed e.g. the following:

- Word level quality analysis for the Finnish part of data
- Open data delivery package of 1771-1910 newspapers and journals (available from digi.kansalliskirjasto.fi)
- Several improvements for the Web interface (time-line, notebook property etc.)
- Ground truth data of Finnish for new optical character recognition (open data)
- A new OCR process with Tesseract 3.04.01
- Named Entity Recognition evaluation collections (two phases: initial trial and present with GT OCR data)
- Layout analysis/article extraction work in progress

# Historical Finnish Newspaper & Journal Web Collection: Digi



THE NATIONAL LIBRARY OF FINLAND

# Three themes of the talk

- Quality of Optical Character Recognition/improvement of OCR
- Named Entity Recognition (NER) on historical Finnish data
- Layout analysis/article extraction work with the data

# Quality: general problems of digitization in old newspaper collections

- Old newspaper collections are hard for digitization (paper and print quality, wear&tear, typeface etc.)

- In the output of the Optical Character Recognition (OCR) process errors are common especially when the texts are printed in the Fraktur (blackletter) typeface.

- E.g. Newspaper collection of British Library has a mean word correctness rate of **about 78 %** (19th Century Newspaper Project, http://www.dlib.org/dlib/july09/munoz/07munoz.html)

- Errors lower usability of corpora both from the point of view of human users as well as regarding potential text mining applications.

# Quality issues of the OCRed Digi data

- Scanning of the contents of Digi was started in the early 2000s

- OCR software etc. was not on the same level then as now

- Quality of the originals has varied quite a lot

- Mostly Fraktur typeface used till the end of 19th century and early 20th century in Finland:

- → quality of the collection is varying due to OCR errors
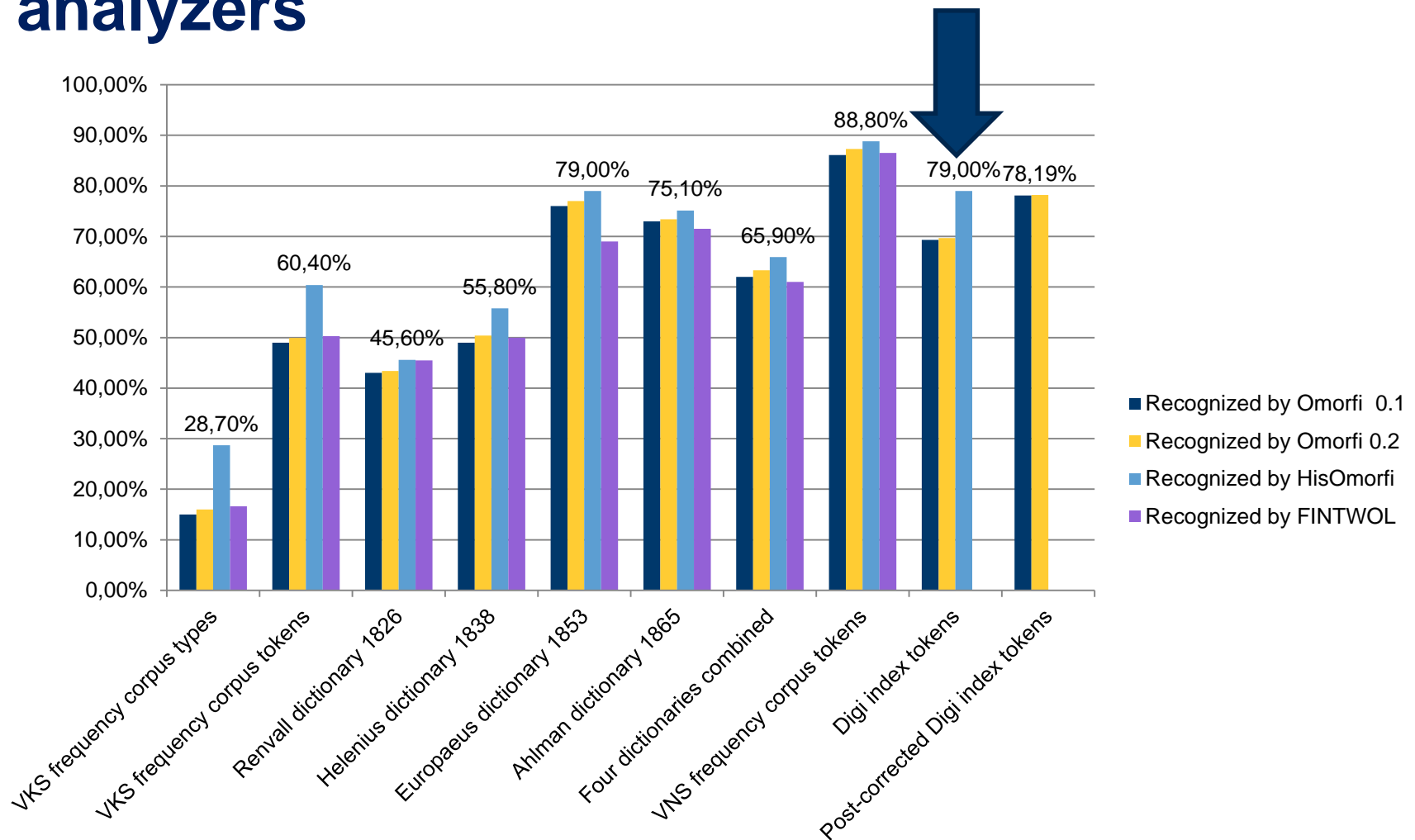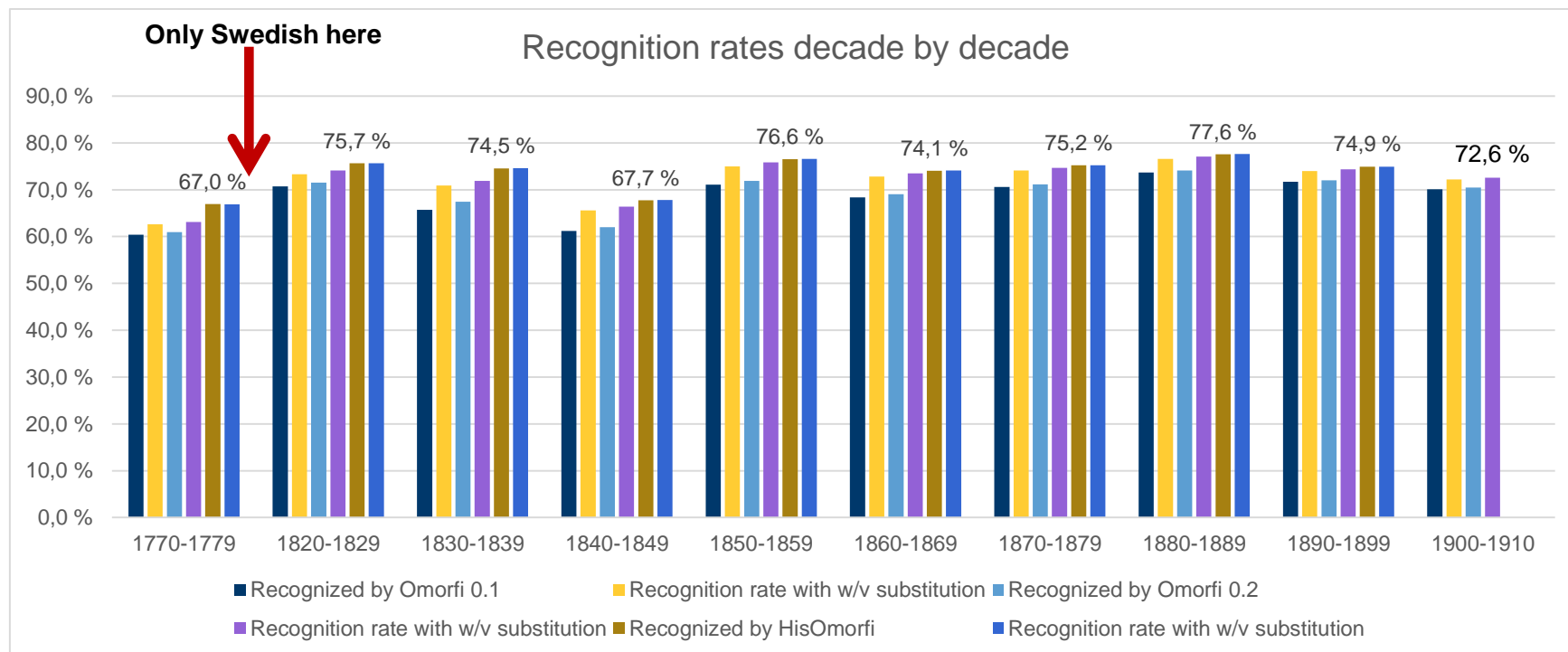
- Aamulehti 17.9. 1905

# Quality assessment

- The Finnish part of the collection 1771-1910 has about 2.40 B words!

- Huge amount - what can you do? (evaluation based on samples or try to assess all?)

- We took four modern Finnish morphological analyzers (FINTWOL and 3 versions of Omorfi) and ran all the word data through them (word data taken out of the web collection's index)

- For comparison we used hand edited comparable word data available from the Institute for the Languages of Finland (VKS_korpus, VNS_korpus and four dictionaries from the 19th century: data is smallish, but the only available: largest c. 4.86 M, dictionaries 25K-90 K)

# Word recognition rates: 4 morphologigal analyzers

# Recognition of newspaper data decade by decade (word tokens)



Recognition rates decade by decade

**Only Swedish here**

Legend:
- Recognized by Omorfi 0.1
- Recognition rate with w/v substitution
- Recognized by Omorfi 0.2
- Recognition rate with w/v substitution
- Recognized by HisOmorfi
- Recognition rate with w/v substitution

Decades: 1770-1779 (67,0 %), 1820-1829 (75,7 %), 1830-1839 (74,5 %), 1840-1849 (67,7 %), 1850-1859 (76,6 %), 1860-1869 (74,1 %), 1870-1879 (75,2 %), 1880-1889 (77,6 %), 1890-1899 (74,9 %), 1900-1910 (72,6 %)

# Caveat – recognition ≠ correctness
## misrecognition ≠ uncorrectness

- *mli*         mli Num Roman Nom Sg → probably an OCR error
- *huu*         huu Part

  *tain*         tai N Gen Sg → wrong division to two parts based on hyphenation, should be *huutain* which is

  unrecognizable, although it is a correct form in 19ᵗʰ century Finnish
- *Hei*         He Pron Nom Pl → should be *heidan,* unrecognizable *(heidän* would be correct)*

  *dan*         +?
- *Samoinkuin*  +? → not recognized because written as a compound, OK otherwise
- *Ylöskannetaan* +? → not recognized because written as a compound, OK otherwise

# Ways to improve the quality

1) Re-OCRing → ABBYY FineReader's Fraktur licensing policy/pricing impossible with new AbbyyFinereader, we are moving to open source Tesseract OCR.

2) Post-correction with software: FIN-CLARIN: about 9-X % units recognition rate improvements so far (Omorfi 0.2, HisOmorfi)

3) Crowdsourcing (i.e. human correction): not feasible due to large amount of data
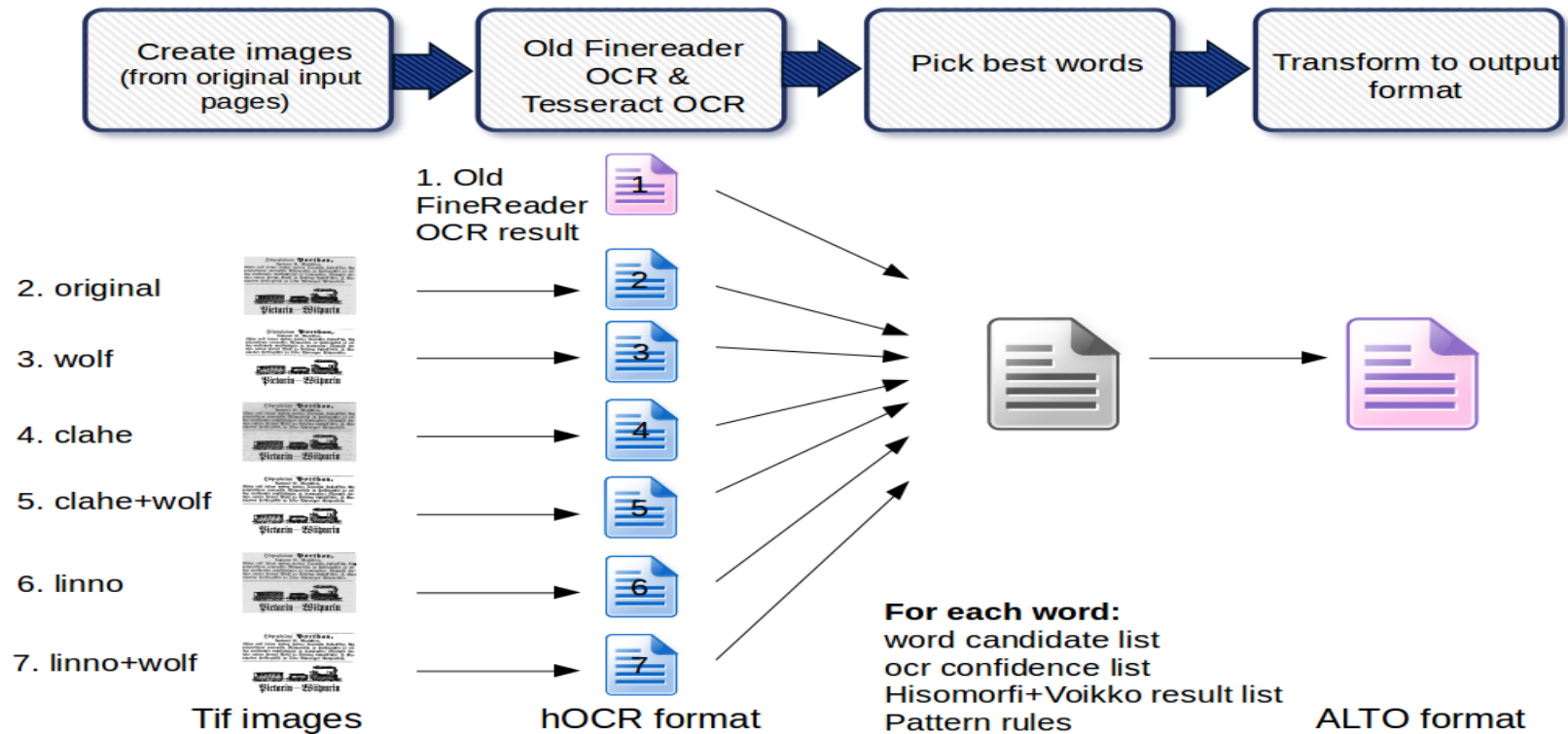
# The new OCR process

- A new Fraktur font model was taught for Tesseract (using an existing German Fraktur as basis)
- A pipeline of image improving process was established
- OCR with Tesseract 3.04.01 (+ existing Abbyy FineReader v.7/8 results)
- Choosing the most probable word candidate from several suggestions
- Transformation of the output to ALTO XML

# The new OCR process

# Results - a 500 K ground truth word list – precision and recall results (latest)

Words without errors 374299
Words with errors 131008
Errorless not corrected 366043
Sum (lines 1 and 2) = 505307
True pos 99071
False negs  31937
False positives =  8256

Recall   =  0.76
Precision =  0.92
F measure=  0.83

Correction rate =  0.69

# Results of 500K GT: character and word error rates, character accuracy

|  | re-OCR of NLF | Current OCR |
|---|---|---|
| CER | 2,05 | 6,47 |
| WER | 6,56 | 25,30 |
| WER (order independent) | 5,51 | 23,41 |
| CAR | 97,64 | 92,62 |

# Uusi Suometar 1869-1890: average gain 14,7 % units in word recognizability

| Year | Current OCR | Re-OCR 2 | Increase in word recognition |
|------|------------|----------|------------------------------|
| 1869 | 69,4 % | 86,59 % | 17,18 % |
| 1870 | 66,98 % | 85,54 % | 18,57 % |
| 1871 | 72,81 % | 87,27 % | 14,46 % |
| 1872 | 75,09 % | 88,25 % | 13,16 % |
| 1873 | 74,61 % | 87,04 % | 12,43 % |
| 1874 | 72,70 % | 86,09 % | 13,39 % |
| 1875 | 70,62 % | 85,52 % | 14,90 % |
| 1876 | 71,50 % | 85,51 % | 14,01 % |
| 1877 | 72,09 % | 84,79 % | 12,70 % |
| 1878 | 70,78 % | 84,70 % | 13,91 % |
| 1879 | 73,52 % | 86,09 % | 12,57 % |
| 1880 | 70,11 % | 85,85 % | 15,74 % |
| 1881 | 67,98 % | 84,26 % | 16,28 % |
| 1882 | 62,41 % | 82,94 % | 20,53 % |
| 1883 | 70,19 % | 82,17 % | 11,98 % |
| 1884 | 69,60 % | 81,67 % | 12,07 % |
| 1885 | 68,11 % | 82,53 % | 14,42 % |
| 1886 | 68,21 % | 82,12 % | 13,92 % |
| 1887 | 65,25 % | 82,16 % | 16,91 % |
| 1888 | 70,27 % | 82,52 % | 12,25 % |
| 1889 | 65,71 % | 81,41 % | 15,70 % |
| 1890 | 64,69 % | 80,71 % | 16,02 % |
| | | | 14,69 % |
| | | | Average |

# NER – Named Entity Recognition

What?

- Names of persons, locations, organisations are important factual data in texts
- They can be recognized automatically to a reasonable extent (70-90+ %)
- They can be used in information extraction out of the data
- Names of persons and locations are used heavily as keywords in text searches of on-line databases. Many times even 80 % of keywords are names.
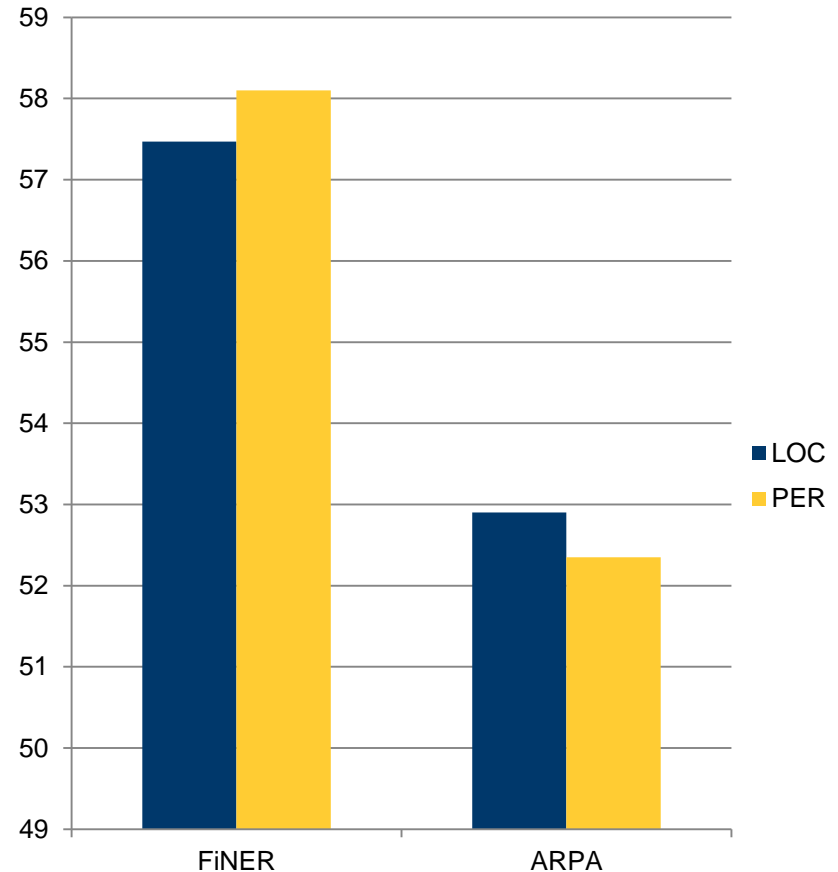
# An example of text with names

Edellä sanotun tarpeen maatimana ja kun ei täällä ole mallikelpoista maanmiljelystaloa, jonka puuttce» poistamiseksi toiwomme että usei»mai»ittu koulu perustettaisiin hra O. L, Vumeruksen omistmnalle Kupiala» tilalle Rantasalmen pitäjään, eikä Iärmikylään Joroisissa, koska Iärmikylä jo kaumcmman aikaa nykyisen omistajansa jalomielisestä ja paljon uhraaivaiscsta toimesta un ollut, ja luonnollisesti tulisi edelleen olemaan PohjoiS-samolaisille monien kokeittcnsa ja mallikelpoisen maanmiljcll,ksen

- 6 names, 2 of them spelled right

# First NER trials with the OCRed newspaper data

- There are tools for doing NER for Finnish, but they are usually for modern language
- We evaluated first 5 different modern tools with an evaluation collection of about 76 000 words (manually tagged by T. Ruokolainen)
- Results were not very

good: F scores 0.52-0.58

(scale 0.0-1.0)

# A new start with Stanford NER package

- We created a new training and evaluation collection of ca. 450 000 words

- Used trainable Stanford NER machine learning package to learn names from about 380 000 words that had been marked manually and semi-automatically

- Results quite good

# Results of Stanford NER

- Ideal results with manually corrected data

- Realistic results with achievable OCR: a 9-10 % unit drop

TABLE I. PRECISION, RECALL, AND F-SCORES FOR EACH NAMED ENTITY CLASS ON THE GROUND TRUTH EVALUATION SET

| Class | Precision | Recall | F1 | # found | #gold standard |
|-------|-----------|--------|--------|---------|----------------|
| LOC | 0.8872 | 0.8566 | 0.8716 | 1764 | 1826 |
| PER | 0.8408 | 0.7801 | 0.8093 | 1118 | 1205 |

TABLE II. PRECISION, RECALL, AND F-SCORES FOR EACH NAMED ENTITY CLASS ON THE OCR EVALUATION SET

| class | precision | recall | F1 | # found | #gold standard |
|-------|-----------|--------|--------|---------|----------------|
| LOC | 0.8527 | 0.7322 | 0.7879 | 1485 | 1826 |
| PER | 0.7856 | 0.6631 | 0.7192 | 1017 | 1205 |

# Results of a LSTM-CRF (Lample et al. 2016)

- LSTM (long short-term memory), recurrent neural network model: state-of-the-art, receives very similar results with Stanford

| Class | Precision | Recall | F1 | # found | # gold standard |
|-------|-----------|--------|--------|---------|-----------------|
| LOC | 0.8598 | 0.6884 | 0.7646 | 1471 | 1826 |
| PER | 0.8212 | 0.6822 | 0.7452 | 1022 | 1205 |

# What to do with NER in a historical newspaper collection?

- NER is a tool that needs to be used for some purpose
- Two basic uses:

  - Enhancement of browsing in large collections
  - Linking names to knowledge sources (DBPedia, Wikipedia, bibliographical sources, geographical sources etc.)

# An example of browsing enhancement

- La Stampa 1867-2005
  http://www.archiviolastampa.it/ww.archiviolastampa.it

# Layout analysis / article extraction

- Newspapers are usually digitised page by page (scanner takes a "photo" of the page)
- Page images serve as the basic unit of whole processing
- → Trouble ahead!

- Page is not any kind of informational unit, only a typographical/printing unit

- Pages consist of different parts: text (news items, titles), pictures etc.

# Uusi Suometar, 16.01.1892, page 4

# Uusi Suometar 31.01.1892, page 4

# Problems

- Number of columns may keep changing even in a single number (and definitely every few years)

- Layout varies: advertisements and other pictures keep popping up in the middle of the text

- Different style of column structure in different parts of page

- ….

Currently best layout analysis software get about 85+% of the layout right in complex layouts/historical data (ICDAR competition from several years)
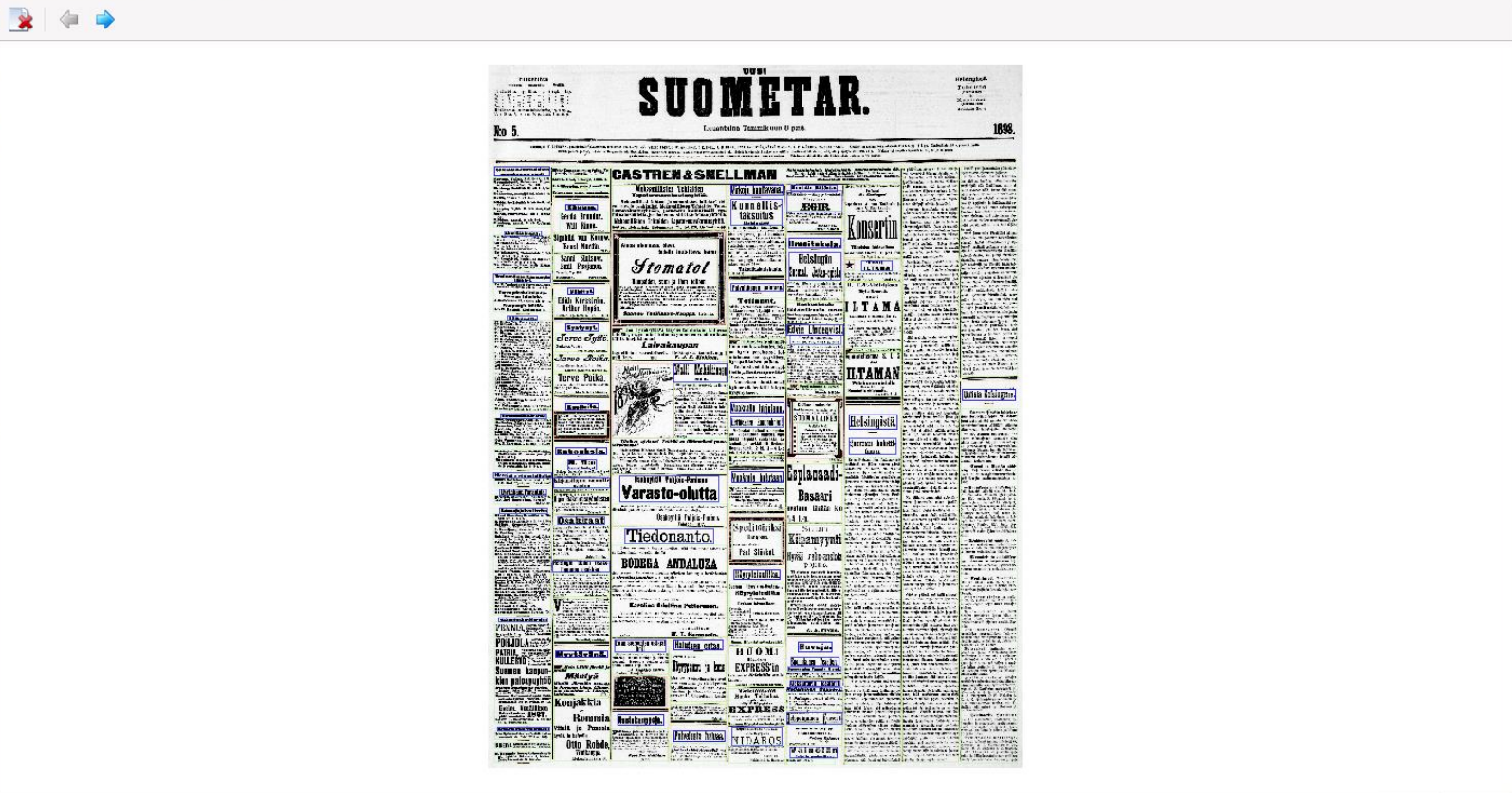
# Article extraction trials at the NLF

- We are collaborating with University of Rouen, LITIS laboratory

- PIVAJ software: a machine learning software, that learns layouts and article separation from training data

- We have made a 224 page manually marked collection out of Uusi Suometar for PIVAJ

- We are trying to train PIVAJ to separate articles from pages of Uusi Suometar

- Results still forthcoming…

# A simple layout page marked with PIVAJ's annotation tool

# Uusi Suometar: different parts marked on a page with PIVAJ in the training data

# In conclusion

- NLF has multilingual newspaper and journal data from 1771-1929 available on the web: digi.kansallikirjasto.fi

- We have been working with data and usability improvement since about 2014 in different projects

- NLF is part of the NewsEye consortium providing data and looking forward to have useful results out of the project

# Thank you for your patience

Kimmo.kettunen@helsinki.fi