# Digitisation and Digital Library Presentation System – A Resource-Conscientious Approach
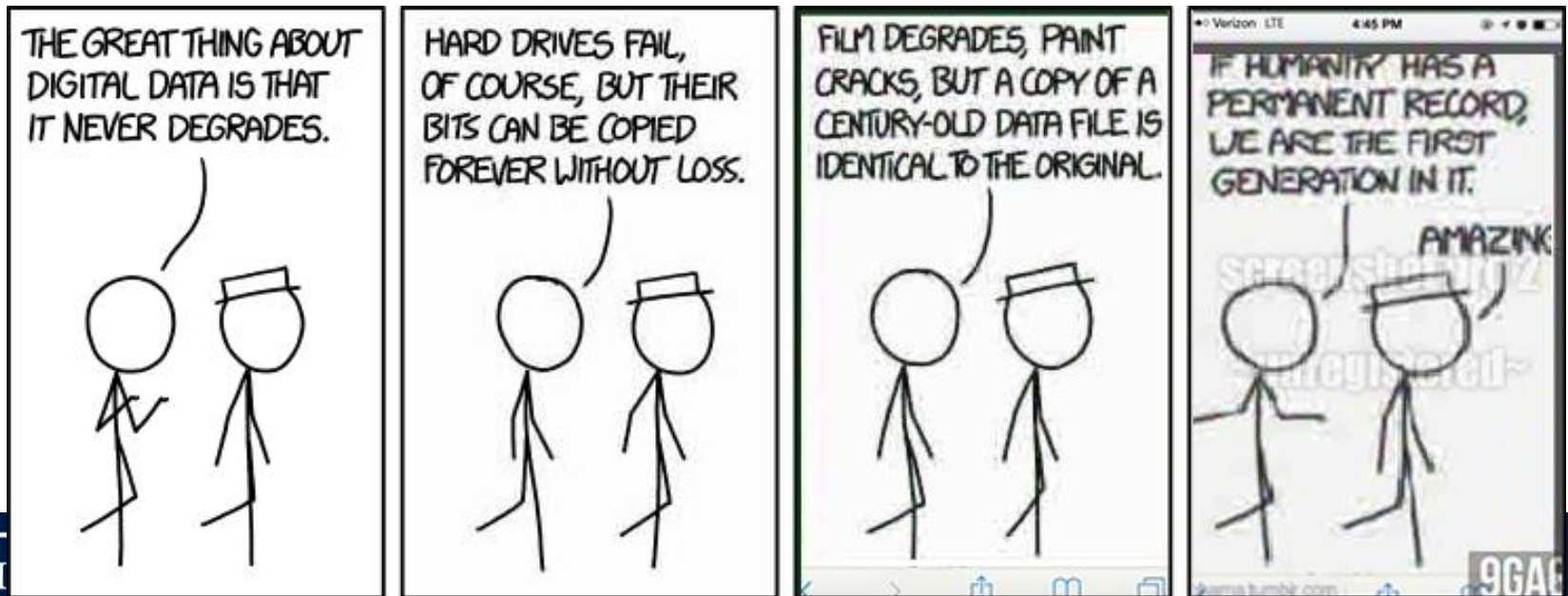
Tuula Pääkkönen, Kimmo Kettunen, Jukka Kervinen,

National Library of Finland

DHN18, 7.-9.3.2018

THE NATIONAL LIBRARY OF FINLAND

# Mission

- Most materials, most easily available  for most users
  - Researchers, citizens, institutions

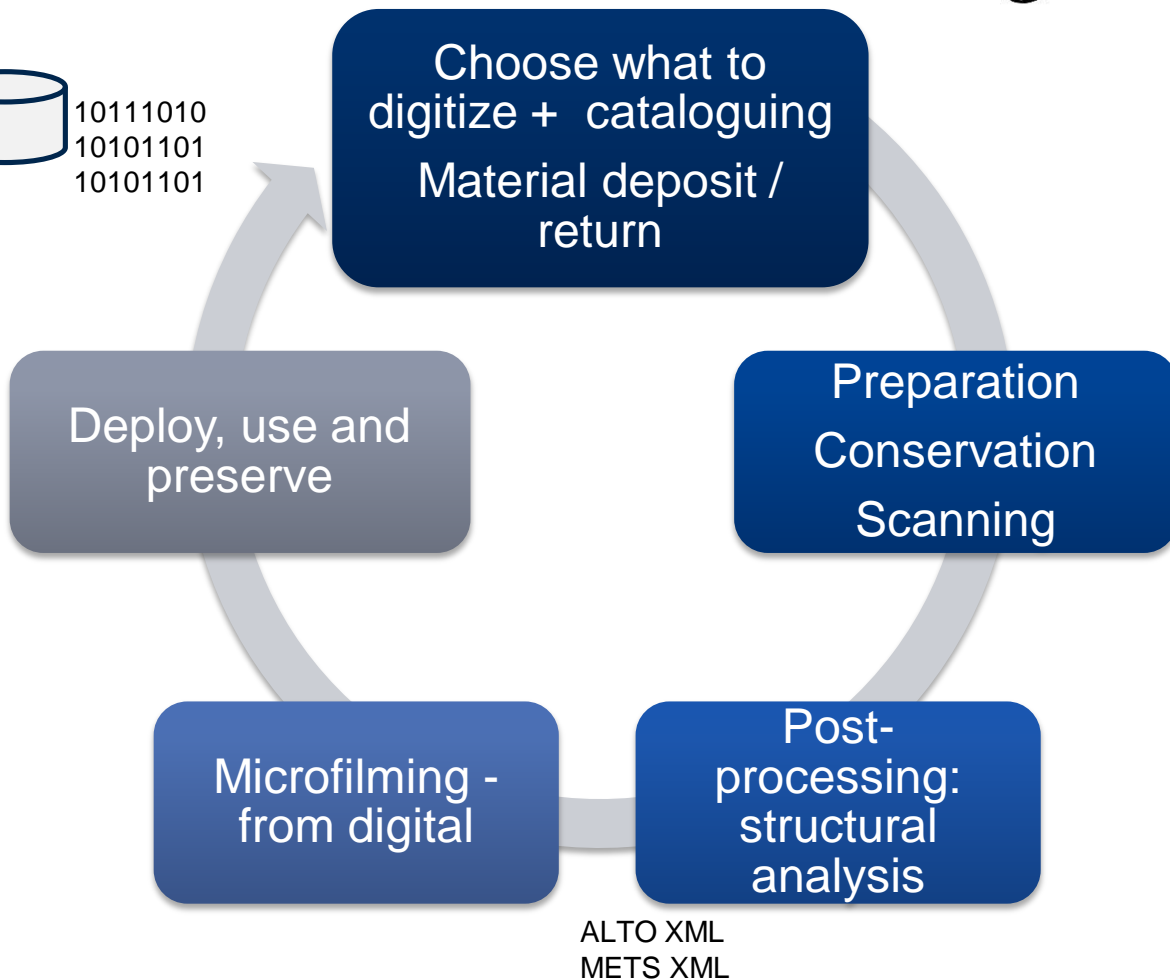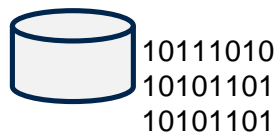- Digitized?    Digitisation policy   &   Time, Quality, Cost



Source: https://xkcd.com/1683/

# Glimpse to Nordic digitization

- Yesterday, we went past 13 million pages digitized (journals and newspapers combined)

| Region | Pages (estimate, million pages) | Material Type | Reference |
|---|---|---|---|
| **Denmark** | 33,3 | Newspapers | http://www2.statsbiblioteket.dk/mediestream/avis |
| **Finland** | 6,0 | Newspapers | https://digi.kansalliskirjasto.fi/info?language=en |
| **Iceland** | 5,4 | Newspapers, periodicals | http://timarit.is/about_init.jsp?lang=en |
| **Norway** | 22 | Newspapers | R. Jøsevold, 'Digitization of copyright protected newspapers in Norway' (Liber 2015) |
| **Sweden** | 17 | Newspapers | http://feedback.tidningar.kb.se/viewtopic.php?id=89 |

# Digital chain



10111010
10101101
10101101

**Choose what to digitize + cataloguing**
**Material deposit / return**

**Preparation**
**Conservation**
**Scanning**

**Post-processing: structural analysis**

ALTO XML
METS XML

**Microfilming - from digital**

**Deploy, use and preserve**

# Digitisation & presentation system together



**Exporter**

**Importer**

**Digitisation**

Export package (zip)
- Preservation image (tiff)
- Access image (jpg, 300dpi)
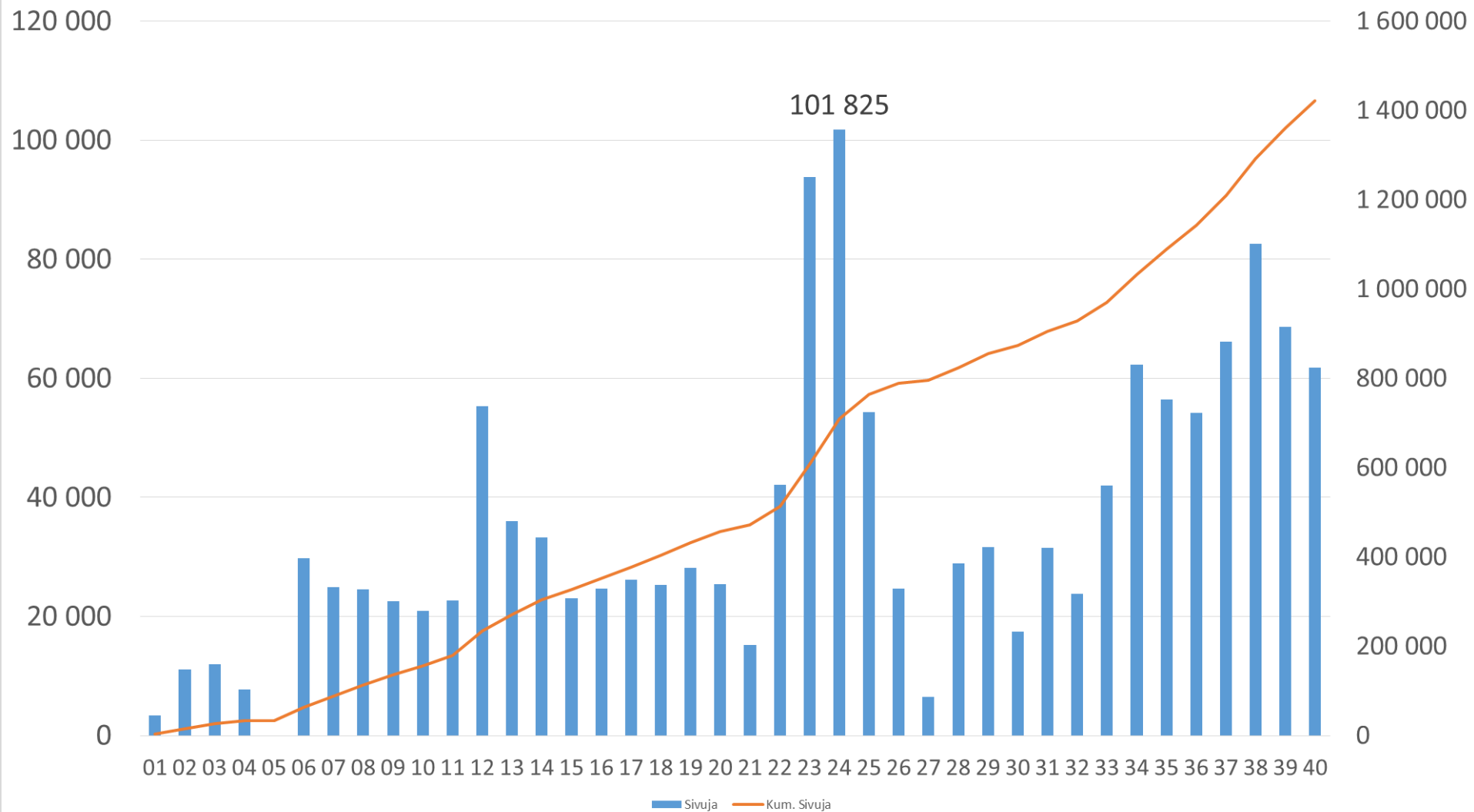- ALTO&METS XML
- Thumbnails
- Pdf

# Audit trail of digitization

- Metadata of post-processing stored



```
1457-4683_1775-09_0_mets.xml ✕        ≡ Untitled-4 ●        ≡ Untitled-1 ●                      ◈   ⊓  ···

335    <mix:ImageCreation>
336        <mix:SourceType>Newspaper</mix:SourceType>
337        <mix:SourceID>MF59418</mix:SourceID>
338        <mix:ImageProducer>The National Library of Finland</mix:ImageProd
339        <mix:Host>
340            <mix:HostComputer>DW19</mix:HostComputer>
341            <mix:OperatingSystem>windows</mix:OperatingSystem>
342            <mix:OSVersion>5.2 Service Pack 2</mix:OSVersion>
343        </mix:Host>
344        <mix:DeviceSource/>
345        <mix:ScanningSystemCapture>
346            <mix:ScanningSystemHardware>
347                <mix:ScannerManufacturer>nextScan,Inc.</mix:ScannerManufa
348                <mix:ScannerModel>
349                    <mix:ScannerModelName>Eclipse Rollfilm</mix:ScannerMo
350                    <mix:ScannerModelNumber>n/a</mix:ScannerModelNumber>
351                    <mix:ScannerModelSerialNo>448013</mix:ScannerModelSer
```
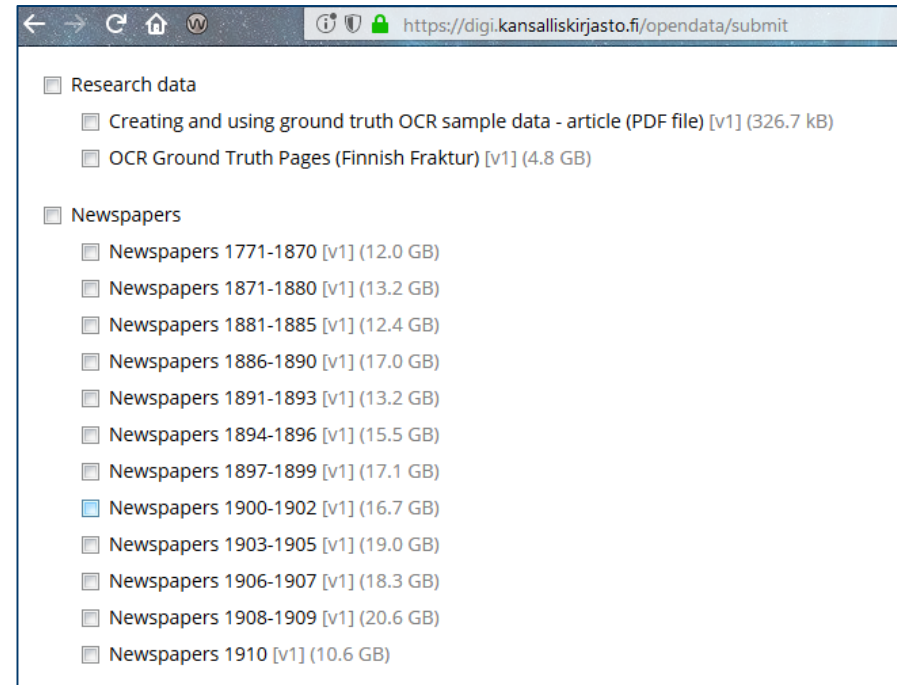
# Pages imported per week



Year 2017 , weeks 1 - 40
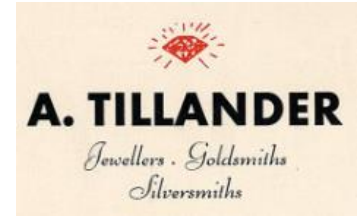
101 825

Sivuja    Kum. Sivuja

# Make room, make room!

- More and more material available as digital/digitized
- New formats and getting existing material digitized requires space somewhere
- Researcher data packages Enrichment work
  - Re-ocr, illustration content analysis, etc.
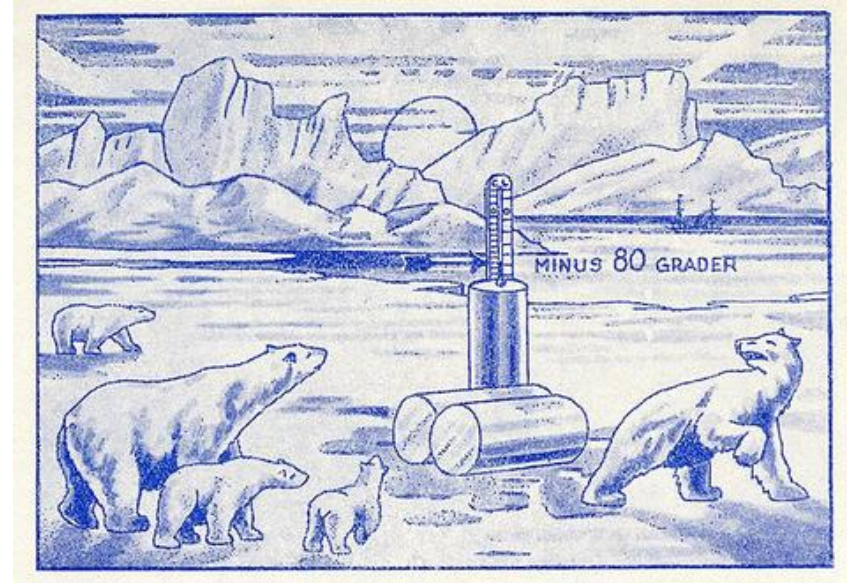
# Final words

- More digitized pages available => more users

- Digitisation is a prioritization issue
  - National
  - Organisational
  - Between material types and projects

- Mass digitization & small-scale digitisation to live together?

# Thank you!

Tuula.paakkonen@helsinki.fi

THE NATIONAL LIBRARY OF FINLAND

# References

- K. Kettunen, E. Mäkelä, T. Ruokolainen, J. Kuokkala, and L. Löfberg, 'Old Content and Modern Tools – Searching Named Entities in a Finnish OCRed Historical Newspaper Collection 1771–1910', *Digit. Humanit. Q.*, vol. 11, no. 3, 2017.

- T. Pääkkönen, J. Kervinen, A. Nivala, K. Kettunen, and E. Mäkelä, 'Exporting Finnish Digitized Historical Newspaper Contents for Offline Use', *D-Lib Mag.*, vol. 22, no. 7/8, Jul. 2016.

- M. Koistinen, K. Kettunen, and T. Pääkkönen, 'Improving Optical Character Recognition of Finnish Historical Newspapers with a Combination of Fraktur & Antiqua Models and Image Preprocessing', in *Proceedings of the 21st Nordic Conference on Computational Linguistics, NoDaLiDa*, Gothenburg, Sweden, 2017, pp. 277–283.

- https://wiki.helsinki.fi/display/Comhis/

THE NATIONAL LIBRARY OF FINLAND