



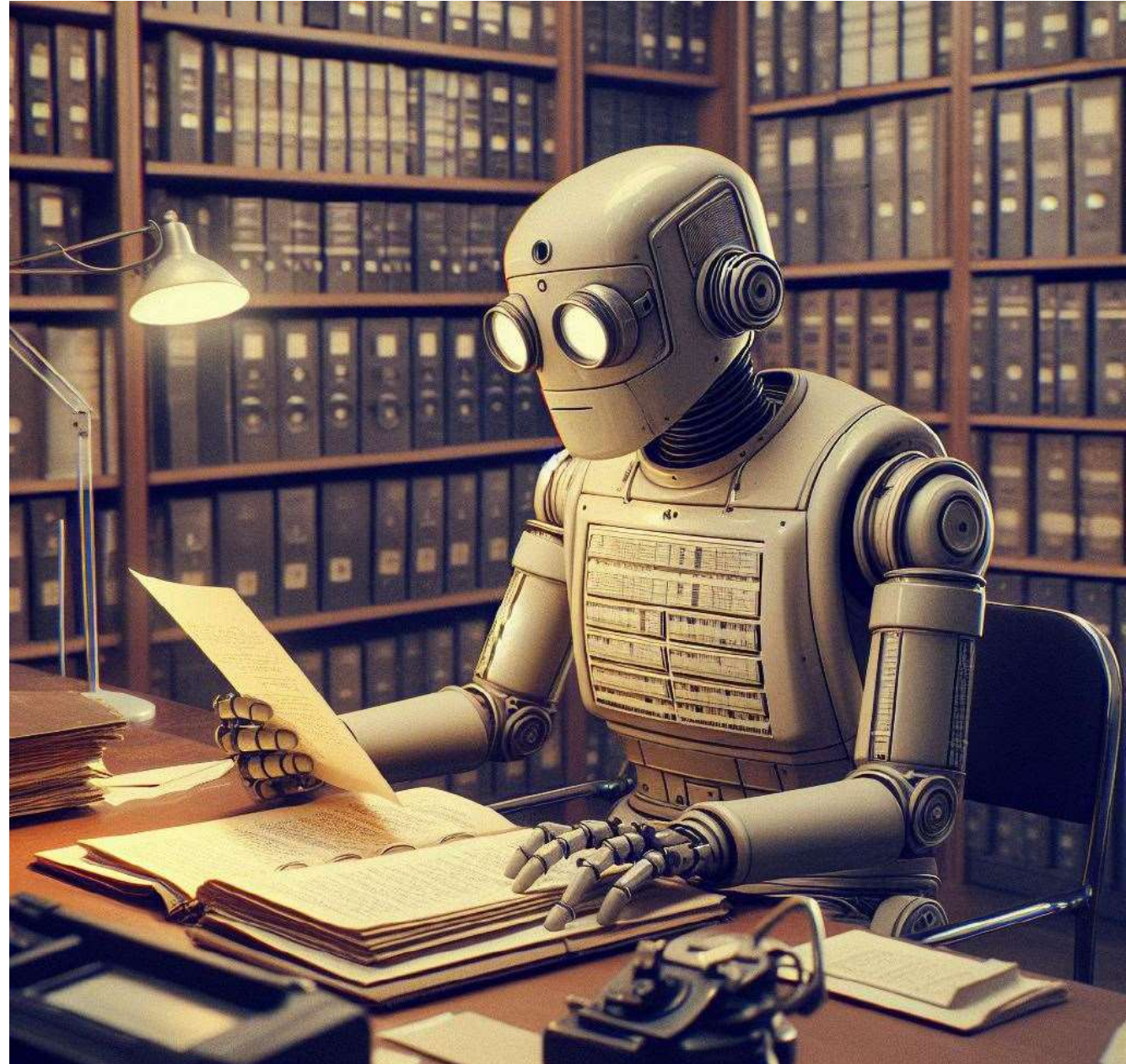
KANSALLISARKISTO

Tekoäly asiakirjahallinnassa

Ilkka Jokipii, yksikönpäällikkö

Mikko Lipsanen, koneoppimisen
pääsuunnittelija

13.4.2024



Luennon sisältö:

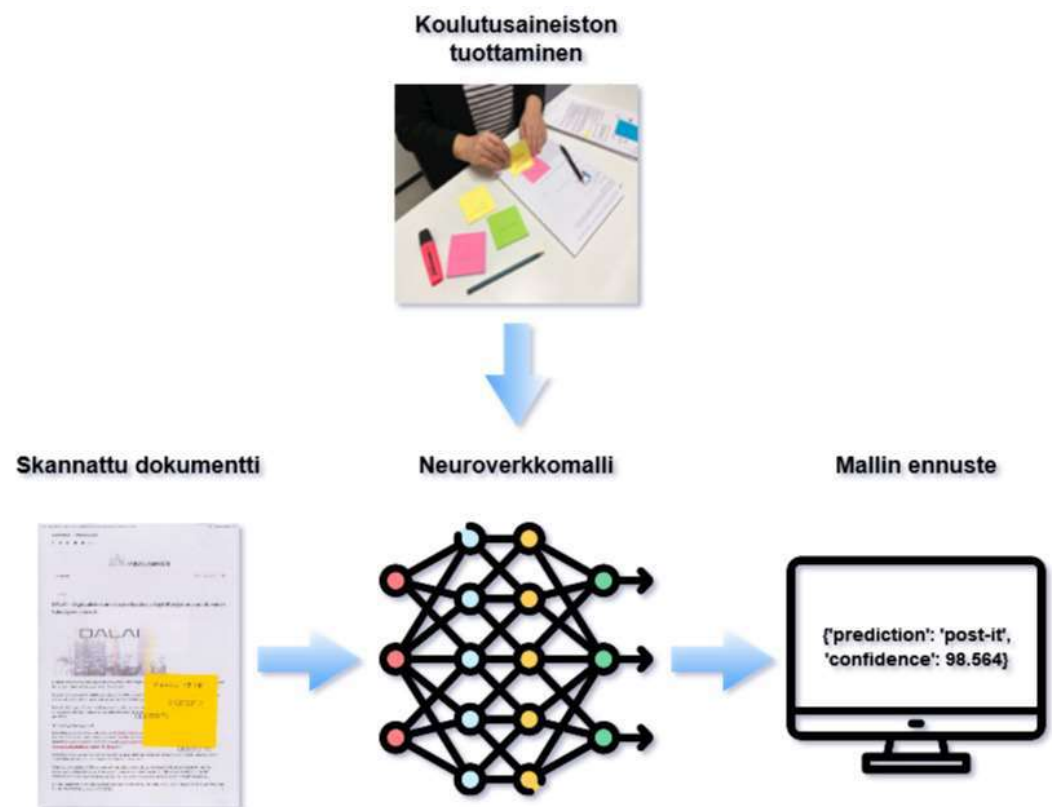
1. Lyhyesti tekoälystä
2. **Tekoälyn hyödyntäminen asiakirjahallinnan eri vaiheissa**
 - **Asiakirjan digitointi**
 - **Digitoidun asiakirjan laadun validointi**
 - **Asiakirjan sisällön tunnistus**
 - **Tunnistetun (teksti)sisällön analysointi**
3. **Käytettävyys: miten tekoälysovelluksiin pääsee käsiksi?**
4. **Tekoäly Kansallisarkistossa**



KANSALLISARKISTO

1. Lyhyesti tekoälystä

- **Tekoäly:** tuo helposti mieleen ihmisen kaltaisen älykkään toimijuuden, vaikka taustalla on yleensä paljon pienimuotoisempia, erilaisiin tehtäviin erikoistuneita koneoppivia sovelluksia
- **Koneoppiminen:** hyödyntää algoritmeja joiden avulla mallit oppivat aineistosta asioita joita niihin ei ole ennalta 'ohjelmoitu': kyky yleistää annetuista lähtökohdista
 - **Ohjattu koneoppiminen:** sovellus oppii etsimään aineistosta haluttuja asioita sille annetun rajallisen esimerkkiaineiston perusteella
 - **Ohjaamaton koneoppiminen:** aineistoa luokitellaan sen piirteiden ja säännönmukaisuuksien perusteella

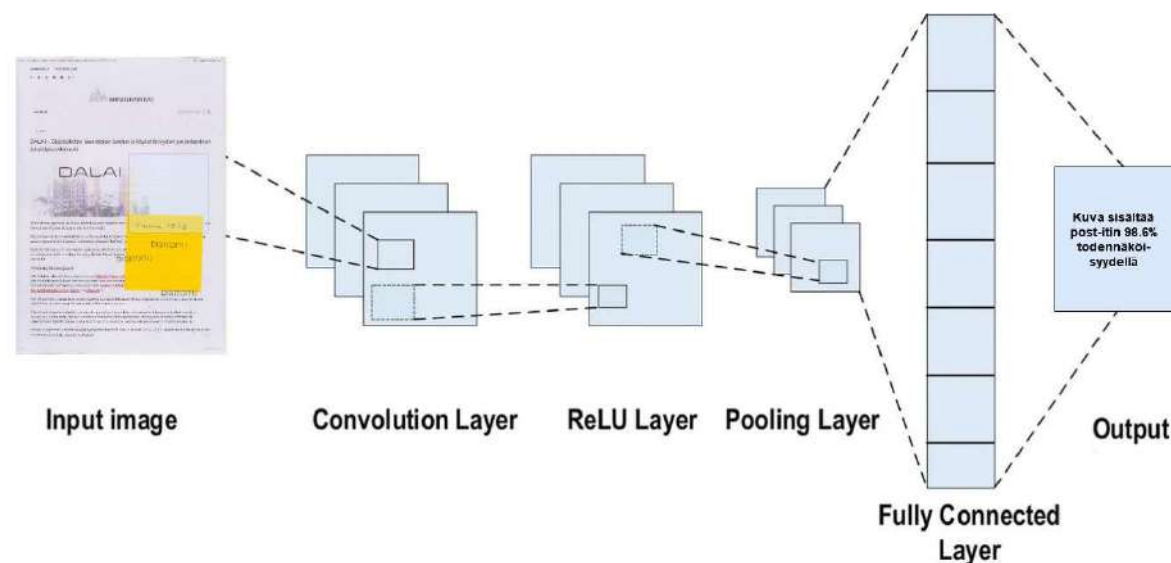


Icons designed by Freepik and xnimrodx from Flaticon

KANSALLISARKISTO

- **Syväoppivat neuroverkkomallit** dominoineet koneoppimisen eri sovelluskohteita (kuvien tunnistus, luonnollisen kielen prosessointi (NLP), robotiikka..) viimeistään 2010-luvun lopulta alkaen

- Tulokset usein ylivoimaisia aiempiin tekniikoihin nähden
- Haasteena kuitenkin mallien vaikea tulkittavuus ('musta laatikko'): päättelylogiikkaa ja virheiden syitä usein vaikea selvittää, kontrollointi ja tulosten tarkkailu tärkeää



Alkuperäinen kuva: Alzubaidi, L., Zhang, J., Humaidi, A.J. *et al.* Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. *J Big Data* 8, 53 (2021).

- **Suurten kielimallien** kehitys tuo mukanaan paljon uusia mahdollisuuksia, mutta niiden hyödyntäminen esim. arkistosektorilla on vielä alkuvaiheessa
 - Voivat helpottaa ja nopeuttaa prosesseja joihin aiemmin tarvittu useita erikseen koulutettuja malleja
 - Toisaalta voivat 'keksiä' omia sisältöjään, kontrollointi usein haastavaa
 - Kaupallisten toimijoiden ja englanninkielisten mallien dominanssi
 - Kehitteillä kuitenkin uusia suomenkielisiä avoimia kielimalleja (TurkuNLP, SiloAI)



Miksi hyödyntää tekoälyä asiakirjahallinnassa?

Tulokset

- Neuroverkkojen kehityksen myötä tekoälymallien nopeus ja tulosten laatu tehneet niistä aiempaa kilpailukykyisempiä.

Yleistvyys

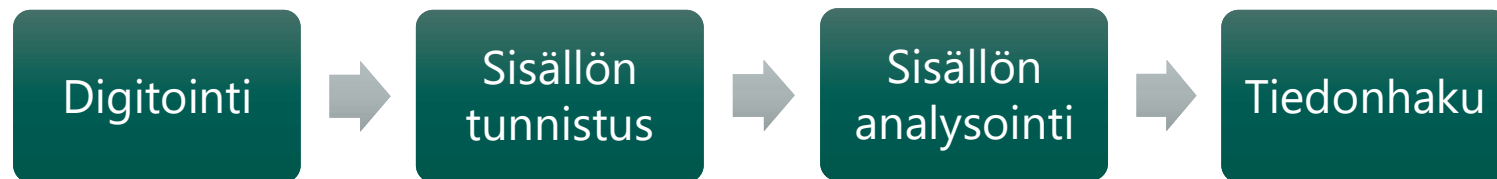
- Manuaalisesti määriteltyihin sääntöihin perustuvia malleja on usein vaikea sovittaa aineiston moninaisuuteen, kun taas koneoppivat mallit pystyvät tähän paremmin.

Käytettävyys

- Koneoppimisen hyödyntäminen on tullut helpommaksi: valmiit sovellukset, esikoulutetut mallit, uudet ohjelmointikirjastot, pilvipalvelut jne.

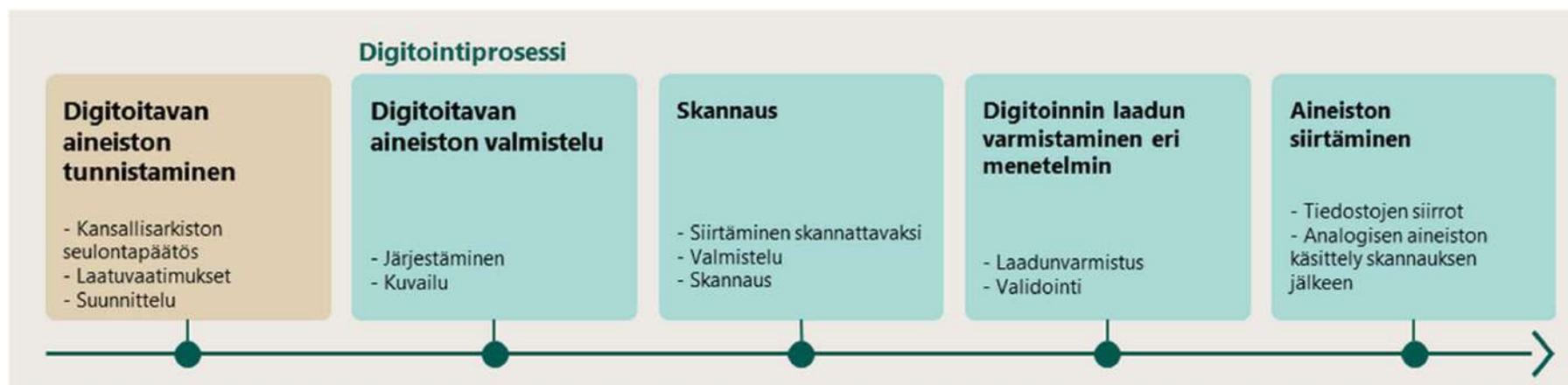
2. Tekoäly asiakirjahallinnan eri vaiheissa

- Käydään läpi tekoälyn hyödyntämisen mahdollisuuksia erityisesti arkistoaineiston näkökulmasta
- Runsaasti käytännön esimerkkejä!



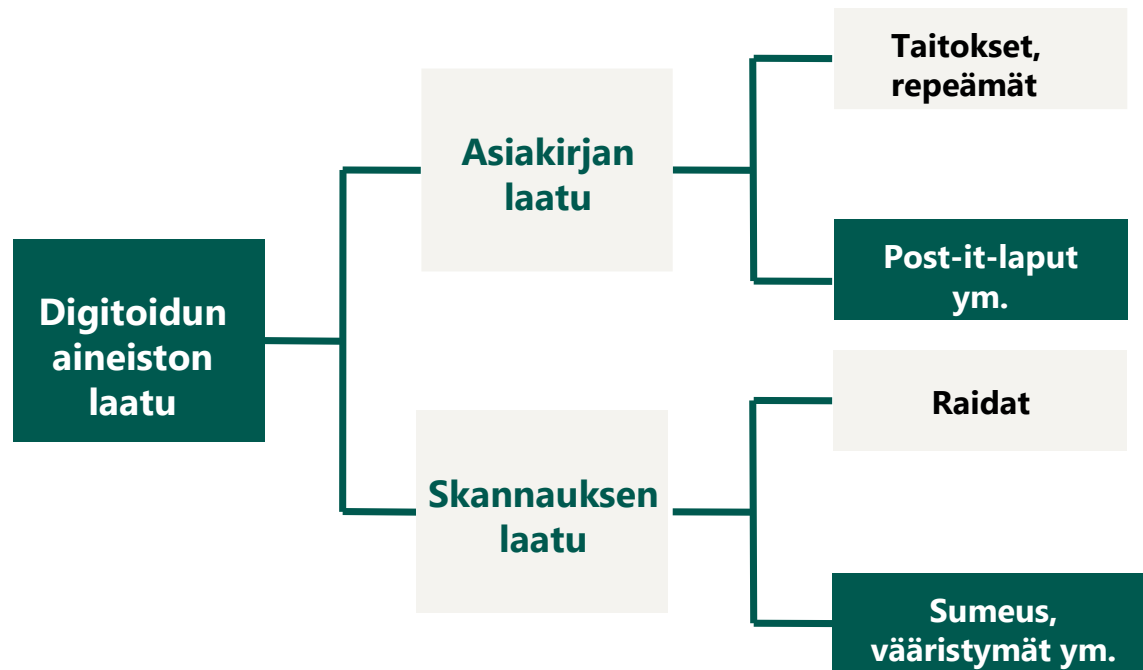
Asiakirjan digitointi

- Digitointiprosessin vaiheet Kansallisarkistossa (<https://kansallisarkisto.fi/digitointi>)



- Tekoälysovelluksilla voi olla rooli jo skannausvaiheessa, mikäli skannerin ohjelmisto hyödyntää sellaisia esim. kuvan laadun tarkkailuun ja optimointiin
- Keskitytään kuitenkin neljänteen vaiheeseen: digitoinnin laadun varmistus ja validointi

Digitoidun asiakirjan laadun validointi



Automatisoinnin hyödyt

- Manuaalinen laadun tarkkailu vaatii paljon resursseja, kun aineistomäärät ovat suuria
 - Esimerkiksi kansallisarkiston massadigitoinnissa* vuosittainen digitoititavoite on 1800 hyllymetriä!
- Automatisoinnin avulla voidaan säästää aikaa ja siten myös rahaa - edellyttäen että käytetyt työkalut toimivat riittävän hyvin

**Massadigitointi: valtion viranomaisten arkistoitavien asiakirjojen digitointi pysyvää säilytystä varten.*

VAIHE 2: Validoidaan aineistoa						Vaihe 1	Määrätyksi	Näytä/hyljät	Leveys 6	Lopetu valkoini
0001	0003	0005	0007	0009	0011					
0013	0015	0017	0019	0021	0023					

Kansallisarkiston massadigitoinnissa käytetty digitoituiden **KANSALLISARKISTO** aineiston manuaalisen validoinnin käyttöliittymä.



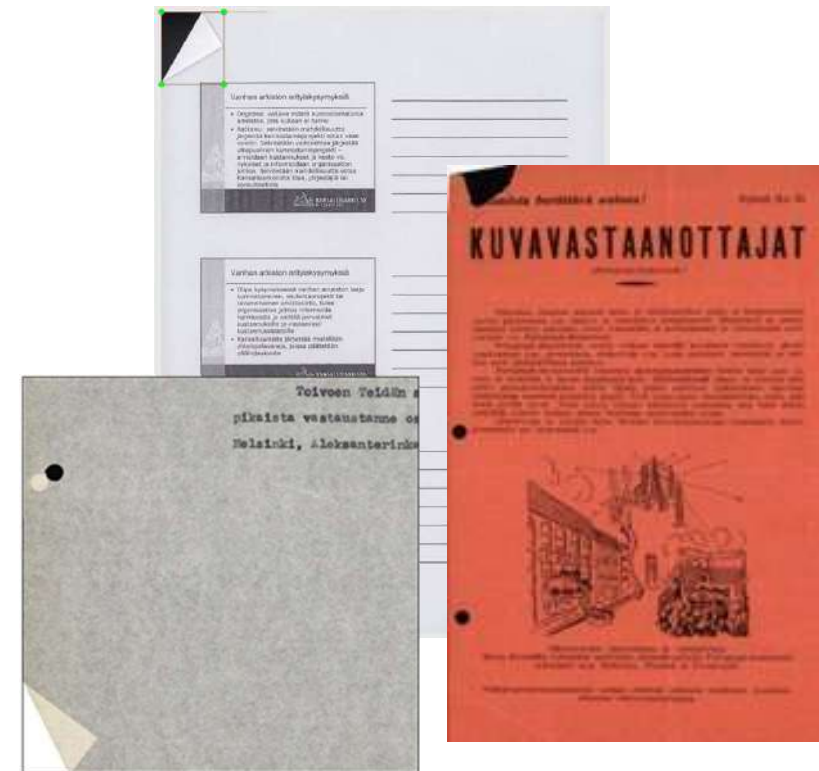
Automatisoinnin haasteet

- Aineiston ja tunnistettavien elementtien moninaisuus on haaste työkalujen yleistettävyydelle
- Koneoppivien mallien koulutukseen tarvitaan oikeanlaista aineistoa ja usein myös manuaalista annotointia*

**Annotointi: aineiston luokittelua / kuvaamista koulutettavan mallin kannalta relevantilla tavalla.*

Taittuneiden kulmien tunnistus

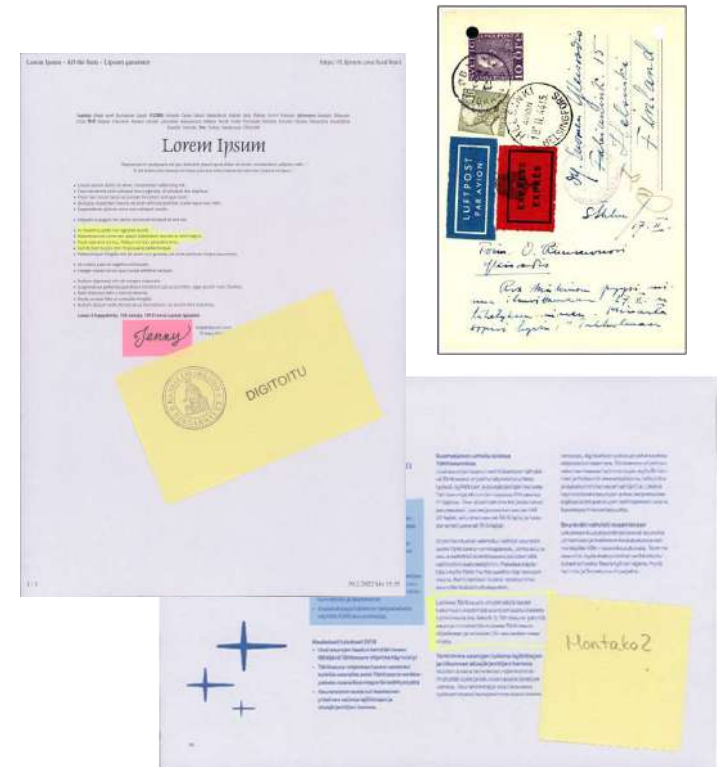
- Sovellus saa syötteenä dokumenttisivun sisältäviä kuvatiedostoja ja luokittelee ne kahteen luokkaan: sivu sisältää taittuneen kulman / sivu ei sisällä taittunutta kulmaa
- Tunnistaa myös repeämät paperin reunoissa
- Neuroverkkomallin koulutusta varten koottiin aineisto jossa oli tuhansia eri tavoin taittuneita ja repeytyneitä asiakirjoja, sekä vastaavasti erilaisia ehjiä asiakirjoja



KANSALLISARKISTO

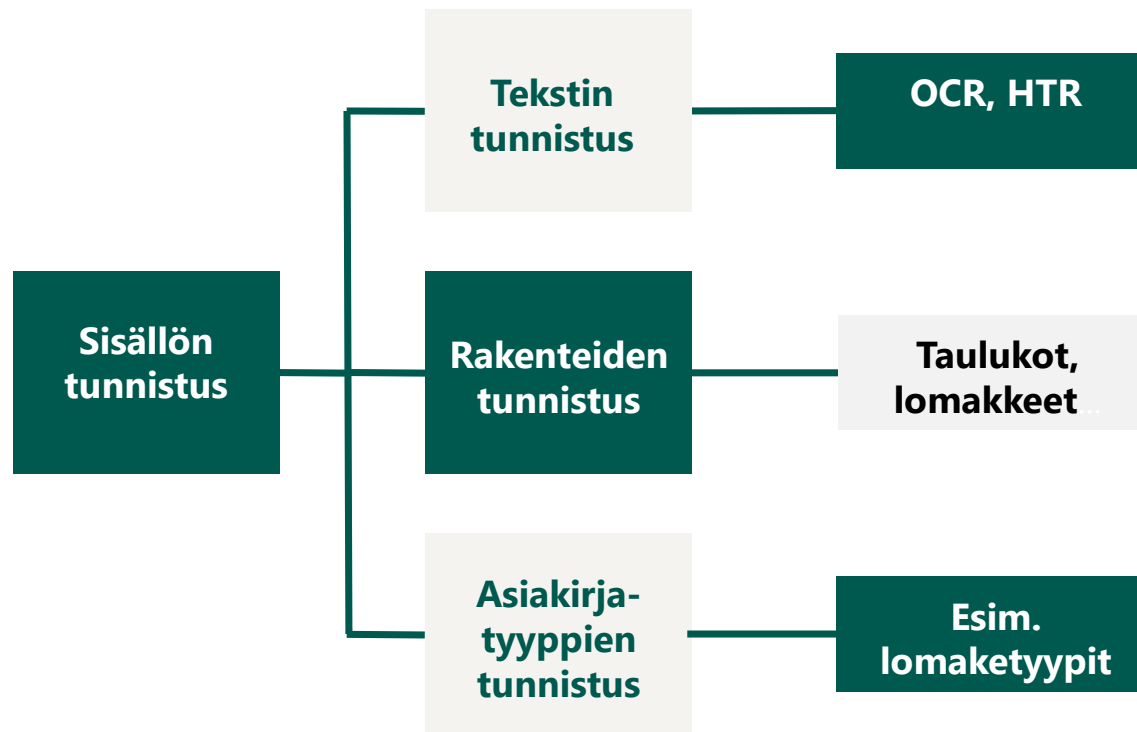
Post-it-lappujen tunnistus

- Digitoija ei aina huomaa post-it-lappuja, jotka voivat peittää dokumentin sisältöä
- Kun post-it huomataan, Kansallisarkiston massadigitoinnissa dokumentti digitoidaan sekä lapun kanssa että ilman sitä
- Sovellus luokittelee sivut kahteen luokkaan: sivu sisältää post-it-lapun / sivu ei sisällä post-it-lappua
- Neuroverkkomallin koulutusaineistoa askarreltiin paljon käsin: eri värisiä ja muotoisia post-it-lappuja tekstillä tai ilman erilaisiin dokumentteihin liimattuina



KANSALLISARKISTO

Digitoidun asiakirjan sisällön tunnistus





Automatisoinnin hyödyt

- Etsityn tiedon sekä aineistojen välisten yhteyksien löytäminen on vaikeaa jos toimitaan vain digitoitujen kuvien tasolla
- Tekstisisällön manuaalinen transkribointi on hidasta: koneoppivat sovellukset nopeuttavat prosessia



Automatisoinnin haasteet

- Tekstisisällön tunnistuksen laatu vaihtelee, vaikuttaa metatietojen tunnistuksen laatuun
- Rakenteelliset elementit kuten taulukot sisältävät valtavasti variaatiota: tunnistuksessa usein toisessa vaakakupissa laatu, toisessa yleistettävyys
 - Tulokset yleensä parempia kun mallit räätälöidään tietyn tyyppiselle aineistolle

Tyhjien sivujen tunnistus

- Lähes 40% Kansallisarkiston massadigitoinnissa digitoiduista sivuista on tyhjiä
- Tyhjiäkään sivuja ei tuhota, mutta niitä ei haluta ajaa sisällöntunnistuksen läpi tai viedä käyttöliittymään
 - Siksi ennen sisällön tunnistusta on tärkeää erotella tyhjät sisällöllisistä dokumenteista
- Tyhjien tunnistusta varten koulutettiin koneoppimismalli, koulutusaineistossa yli 100 000 kuvatiedostoa
- Jopa tyhjän sivun tunnistus voi olla yllättävän haastava tehtävä: esimerkiksi sivun läpi heijastuva arkin kääntöpuolen teksti aiheuttaa ongelmia



Konekirjoitetun tekstin tunnistus

- **OCR** (Optical Character Recognition) on keskeinen vaihe digitoidun asiakirjan prosessoinnissa
 - Tekstisisällön muuttaminen digitaaliseen muotoon manuaalisesti erittäin aikaa vievä prosessi
 - OCR:n laatu vaikuttaa kaikkiin prosessoinnin myöhempiin vaiheisiin, joissa hyödynnetään asiakirjan tekstisisältöä
- Tarjolla lukuisia ilmaisia (Tesseract, EasyOCR, PaddleOCR..) ja kaupallisia (Amazon Textract, Microsoft Read OCR, Google Document AI..) sovelluksia
- Mitä rakenteisempi asiakirja on, ja mitä heikompi asiakirjan tai skannauksen laatu, sitä huonolaatuisempi on myös OCR:n tulos
- Tesseract yleisesti (myös Kansallisarkistolla) käytetty ilmaissovellus
- Kansallisarkisto mukana hankkeessa jossa tutkitaan, voidaanko sovellusten laatua parantaa jatkokouluttamalla niitä omalla aineistolla

Konekirjoitetun tekstin tunnistus

Lomake
Blankett No 401

Posti- ja lennätinlaitos — Post- och telegrafverket

Mistä *)
Från *) LONDON LH134 36 4 1015 NORTHERN =

Virkamultatukata
Tjänstemärkningar

Sähkösanoma — Telegram	Vastaanotettu — Emottaget
ELT = MADAME WUOLIJOKI STATSRADION HLS = 232478	1 10 kl. t. m. -4 II. 1947 192 Johd. No ledn. No

Numerot lähetyspaikan nimen jälkeen merkitsevät: 1) sähkösanoman numero, 2) sanaluku, 3), 4) ja 5) sähkösanoman sisältöjätövä, -tuntia ja -minuuttia.

*) Sifferna efter avsändningsorten betyder: 1) telegrammets nummer, 2) ordantal, 3), 4) och 5) inlämningsdatum, -timme och -minut.

SUGGEST YOU HEAR ELIZABETH LOCKHART DANISH
RADIO 9=TH 2110 DANISH TIME STOP SHE'S
RECITING STOCKHOLM 20=TH STOP ROBERT HAS
REPERTOIRE IF INTERESTED ENGAGING HER RADIO
PLUS OTHER ENGAGEMENTS JUSTIFYING BALTIC
CROSSING = ERIC DANCY

*Antarut
hitid samant*

*British Council ei järjestä tähän otteen
vaan myöhemmin.*



Lomake
IFA No 401

Posti- ja lennätinlaitos — Post- och telegrafverket

Mistä *)
Från *) EL LONDON LH134 36 4 1015 NORTHERN c-
Virkamultatukata I
Tjänstemärkningar A

Sähkösanoma — Telegram Vastaanotettu — Emottaget

ELT = MADAME WUOLIJOKI k sseenetd [p
E E VO A 5/9. 6 a
9 - omitt A HMLUCV TY T 11 ' (- e, k:
ta SLATSRADION HLS - i H. 1947 192 k
3N | - E
2:85 VY * III 7 18 C, p E
4 PILP 4f JOHT.
JBS ad ledn. VO) E

3N lähetyspaikan nimen jälkeen merkitsevät: 1) sähkösanoman numero, 2) sanaluku, 3), 4) ja 5) sähkösanoman sisältöjätövä, -tuntia ja -minuuttia.

*) Sifferna efter avsändningsorten betyder: 1) telegrammets nummer, 2) ordantal, 3), 4) och 5) inlämningsdatum, -timme och -minut.

ROT VA 1 11 N I TGA ra "51 KA MITA Pit TIAN FC 1
> US TUB D FOU MII AN L ANDIALIUI U G BIN TA 4 IAIN LO PI
S 000 + 87 | [5 va TY AN TOI MI TIPTG SC'A 12 (6.3 4 X <
RADIO J=1 [1 2110 J AIM LOTI A LK DA TT 1) 4-1 1)
n TTI TEE m / NA 29=1 H STOP m RS mu) (1
IV CLIN 6 STOCKHO Kl Ca = m VOE N TR YU
m - * %w fy (3 J
Ora 3: 1971 n Y ra v AIM MMA PCI ra TY IA 5 tr 1 AT) 7
Kia PRI OIRN IF LIN LIKI D 4 4 J TNGAG ING HMR KADJIO
, 4 TAV TAISEN TV A m T A
(4 YT 1015) ME A ffitah mj M to TII V i PA
U: AALIKSN. TNGAT UNA AN i JUSTIEF Y 11C- BALTIC
D (6 C TNT DT TIAN (VY
KUDU ANT E IT RIU JAN J

/ -hrandermdg
ae / Pe istot Powel

V. e na att a
/ R >> f (punti EA Vary 7 (n ja ha see dt
) Kin 4 h (')

KANSALLISARKISTO

Muiden sisältöelementtien tunnistus

- Neuroverkkoihin pohjautuvat objektintunnistusmallit mahdollistavat monen tyyppisten elementtien tunnistamisen asiakirjoista
- Valokuvat, piirroksot, kartat, leimat, taulukot, allekirjoitukset...
- Luokittelu: asiakirja sisältää kuvan
- Paikannus: kuva sijaitsee tietyssä kohdassa dokumenttia
- Esimerkiksi asiakirjat sisältävät usein sekä käsin- että konekirjoitettua tekstiä, joiden tulkintaan käytetään omia koneoppimismalleja
 - Kansallisarkistolla testattu käsin- ja konekirjoitettujen tekstikohtien tunnistusta ja erottelemista dokumentista kohdennettua tekstintunnistusta varten



KANSALLISARKISTO

Asiakirjan rakenteiden tunnistus

- Laadukas tekstisisällön tunnistus edellyttää usein rakenteiden tunnistusta: esim. taulukon sisältävän dokumentin syöttäminen suoraan ocr-sovellukselle tuottaa yleensä varsin sekavia tuloksia
- Toisaalta dokumentin rakenteet kuten taulukot ja lomakepohjat halutaan usein muuntaa digitaaliseen muotoon mahdollisimman uskollisesti: tällöin täytyy tunnistaa rakenne ja teksti sekä sijoittaa teksti oikein rakenteeseen
- Rakenteiden tunnistukseen tarjolla erilaisia tekniikoita ja sovelluksia, tässäkin tapauksessa ei ole olemassa yhtä tiettyä dominoivaa menetelmää

The image shows a complex form with multiple sections and tables. The top section is titled 'PERINTÖVEROKUUKSEN VALMISTELULOMAKE'. Below it, there are several tables and form fields. The tables contain numerical data and text. The form fields are for text entry. The overall layout is dense and organized into distinct sections.

The image shows a large, multi-column table with a grid structure. The columns are labeled 'Henkilöluokka' and 'Perustiedot'. The table contains a large amount of text and numbers. The overall layout is dense and organized into distinct sections.

Tunnistetun (teksti)sisällön analysointi

- Aineiston 'raakasisällön' mahdollisimman laadukas tunnistus on keskeinen vaihe asiakirjojen saattamisessa käyttäjille digitaalisessa muodossa
- Yhä tärkeämmäksi on kuitenkin muodostunut myös tunnistetun sisällön analysointi ja jalostaminen erilaisia käyttötarkoituksia varten
- Käyttäjä voi esimerkiksi haluta etsiä aineistosta tiettyjä henkilöitä tai paikkoja, etsiä tiettyyn paikkaan, tapahtumaan tai aihepiiriin liittyviä asiakirjoja, lukea tiivistelmiä aineistokokonaisuuksien sisällöstä, tai vaikkapa esittää sovellukselle kysymyksiä aineistoon liittyen
- Tässä tekoälyn ja erityisesti suurten kielimallien kehitys tarjoaa monenlaisia mahdollisuuksia



KANSALLISARKISTO

Nimientiteettien tunnistus

- Mistä NERissä on kyse?
 - Tunnistetaan tekstidatasta ennalta määriteltyihin kategorioihin kuuluvia elementtejä
 - Yleisiä entiteettejä esim. henkilön nimet, paikannimet, organisaatiot, aikamääreet..
 - Hyödyllinen työkalu digitoidun aineiston automaattisessa metatiedottamisessa

Valtion Lisenssitoimikunnalle **ORG** . Saksalainen **NORP** kirjailija, tri Friedrich Wolf **PERSON** , Berlin **GPE** - **PERSON** Pankow **GPE** , on Suomen Yleisradiolle **ORG** luovuttanut 2 käsikirjoituksensa esitysoikeuden, joista käsikirjoituksista toinen on jo esitetty ohjelmistossa ja toinen tullaan lähiaikoina esittämään. Esitysoikeuden luovuttamisesta on tri Wolf **PERSON** saamassa yhtiöltämme pakkionsa, jonka hän meille saapuneen tiedon mukaan toivoisi saavansa vastaanottaa elintarvikkeiden muodossa. Tämän vuoksi Oy.Yleisradio Ab. **ORG** anoo kunnioittavasti, että Valtion lisenssitoimikunta **ORG** myöntäisi yhtiöllemme vientilisenssin 10 kg:lle lihaa, sen lähettämistä varten tri Wolfille **PERSON** . Lisenssit pyydämme saada kahta 5 kg:n lähetystä varten. Helsingissä **GPE** , tammikuun 20 p:nä 1948 **DATE** . Hella Wuolijoki **PERSON** . Einar Sundström **PERSON** .

- NER-mallit yleensä koulutettu datalla joka
 - a) ei sisällä juurikaan asiakirja-aineistoa ja
 - b) ei sisällä tekstintunnistuksen (OCR, HTR) tuottamaa 'meluista' aineistoa
- Tästä syystä myös tunnistavat heikosti entiteettejä tällaisesta aineistosta
- Kansallisarkisto toteutti yhdessä Fin-Clariah –tutkimuskonsortion tutkijoiden kanssa erityisesti arkistoaineistoa varten kehitetyn nimientiteettien tunnistusmallin
- Tunnistaa yhteensä 10 luokkaan kuuluvia nimientiteettejä:
 - Henkilönnimet, organisaatiot, paikat, geopoliittiset alueet, tuotteet, tapahtumat, päivämäärät, kansallisuudet ja uskonnolliset tai poliittiset ryhmät, diaarinumerot, y-tunnukset
- Yhdistetyllä arkistoaineistolla sekä uudemmalla (mm. uutistekstejä, blogitekstejä, lakitekstejä) datalla koulutettu malli tunnistaa entiteettejä hyvin eri tyyppisistä aineistoista, hälyisestä OCR:äytystä asiakirjatekstistä tai digisyntyisestä blogitekstistä

```

Pelkkää O
tyhjyyttä O

Kävin O
tänään B-DATE
katsomassa O
Suomen B-ORG
Perinteisen I-ORG
Teatterin I-ORG
näytelmän O
Ranta B-WORK_OF_ART
. O

Jo O
teatterin O
nimi O
antiklimaksi O
; O
ikäänkuin O
ryhmä O
tahtoi O
antaa O
katsojalle O
vakuutuksen O
siitä O

```

Esimerkki Turku OntoNotes Entities Corpus-aineistosta

KANSALLISARKISTO

- **Elian Lehto** sijoittui 11:nneksi Kitzbühelin maailmancupin syöksyssä. Yksikään suomalainen mieslaskija ei ole sijoittunut syöksyssä yhtä korkealle maailmancupissa. **Andreas Romar** sijoittui aikoinaan Bormiossa 2012 sijalle 14.



Elian Lehto **PERSON** sijoittui 11:nneksi Kitzbühelin maailmancupin **EVENT** syöksyssä. Yksikään suomalainen **NORP** mieslaskija ei ole sijoittunut syöksyssä yhtä korkealle maailmancupissa. Andreas Romar **PERSON** sijoittui aikoinaan Bormiossa **GPE** 2012 **DATE** sijalle 14.

Herra Wiljo Lindbohm,
Mikkeli.

Vastaukseksi arv. 2/8-31 päivättyyn kirjeeseen saamme
 täten ilmoittaa varanneemme Teille 25 minuuttia syysk. 7 pnä klo
 18.00-18.25.

Kunnioitavasti
 AMJ.



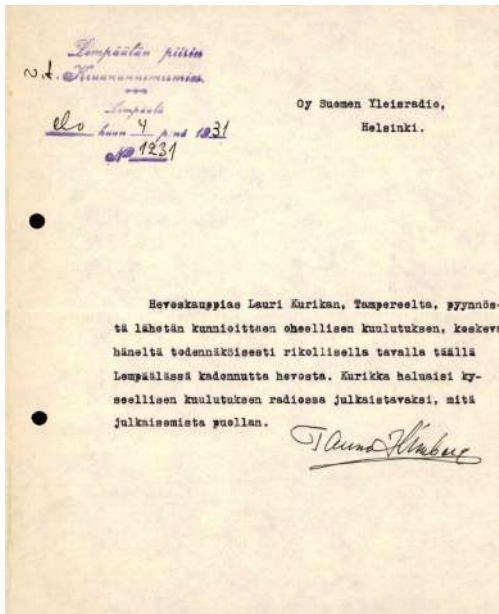
Mikkeli **GPE** . , 13 pni elokuuta = 31 **DATE** ... Herra Wiljo Lindbohm **PERSON** , N ; N j Vastaukseksi arv. 2/8-31 **DATE** päivättyyn kirjeeseen SN S. täten ilmoittaa varanneemme Teille 25 minuuttia syysk. 7 pnä **DATE** klo a W : 18.00+18.25. | : | i 4 N Kunnioitavasti Eeti W An: dä s R

KANSALLISARKISTO

Automaattinen asiasanoitus

- Nimientiteettien tunnistusmalli etsii asiakirjoista vain niitä ennalta määriteltyjä entiteettejä joilla malli on koulutettu
 - Ei ole yleensä paras vaihtoehto aineiston sisällön kuvailuun
- Automaattinen asiasanoitus tuottaa syötteenä saadun tekstin pohjalta sitä kuvailevia asiasanoja
- Tunnettu suomenkielinen sovellus Kansalliskirjaston toimesta kehitetty **Annif** (<https://annif.org/>)
 - Hyödyntää koneoppivia neuroverkkomalleja ja asiasanalistoja (esim. YSO - Yleinen suomalainen ontologia)
 - Annifia käyttää mm. automaattinen sisällönkuvailun palvelu Finto AI (<https://ai.finto.fi>)

Automaattinen asiasanoitus



PA
(W 4 6 €.) 6 04 64 91 4678 906071 6Å 199664,
Oy Suomen Yleisradio,
«i tmÄ Kala
o kuun A A:nd 1931 Helsinki.

Hevoskauppias Lauri Kurikan, Tampereelta, pyynnöstä lähetän kumioitteen oheellisen kuulutuksen, koskeva häneltä todennäköisesti rikollisella tavalla täällä Lempäälässä kadonnutta hevosta. Kurikka haluaisi kyseellisen kuulutuksen radiossa julkaistavaksi, mitä julkaisemista puollan. ET A / v

Get suggestions → **annif**

SUGGESTED SUBJECTS

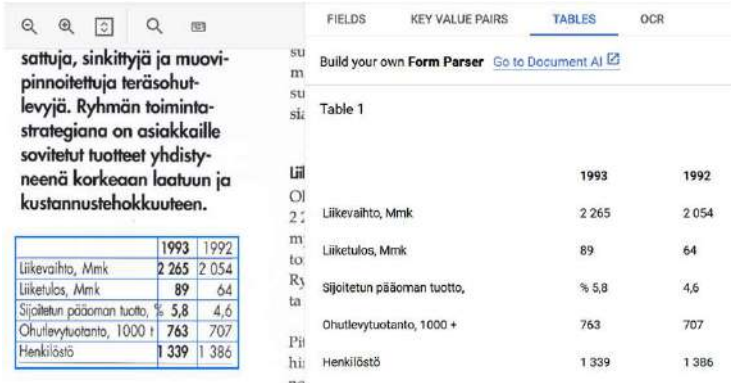
- historia
- Kurikka
- Lempäälä
- suomen kieli
- muistelmat
- yleisradiotoiminta
- radio-ohjelmat
- paikallishistoriat (historiikit)
- hevonen
- paikannimet

Digitoitu asiakirja, tekstintunnistuksen (Tesseract OCR) siitä tunnistama tekstisisältö sekä automaattisen asiasanoituksen (Annif) tekstin pohjalta tuottamat asiasanat.

KANSALLISARKISTO

3. Käytettävyys: miten tekoälysovellukseen pääsee käsiksi?

- Suurilla pilvipalvelun tarjoajilla on omat alustansa joissa saatavilla erilaisia (pääosin maksullisia) sovelluksia asiakirjojen automaattiseen prosessointiin. Esimerkiksi
 - Amazon Textract: kone- ja käsinkirjoitetun tekstin ja dokumentin rakenteiden tunnistus
 - Amazon Comprehend: kielen tunnistus, (tiettyjen) nimientiteettien tunnistus, dokumenttien luokittelu (positiivinen / negatiivinen sisältö)...
 - Microsoft Azure AI Document Intelligence: dokumentin tekstisisällön ja rakenteiden tunnistus...
 - Microsoft Azure AI Language: (tiettyjen) nimientiteettien tunnistus, dokumenttien luokittelu (positiivinen / negatiivinen sisältö), tiivistelmien teko...
 - Google Document AI: kone- ja käsinkirjoitetun tekstin ja dokumentin rakenteiden tunnistus...



The screenshot shows the Amazon Textract interface. On the left, there is a preview of a document with Finnish text: "sattuja, sinkittyjä ja muovipinnoitettuja teräsohuttelevyjä. Ryhmän toimintastrategiana on asiakkaille sovitut tuotteet yhdistyneenä korkeaan laatuun ja kustannustehokkuuteen." Below the text is a table with three columns: the first column lists categories, and the next two columns show values for the years 1993 and 1992.

	1993	1992
Liikevaihto, Mmk	2 265	2 054
Liiketulos, Mmk	89	64
Sijoitetun pääoman tuotto, %	5,8	4,6
Ohutlevytuotanto, 1000 +	763	707
Henkilöstö	1 339	1 386

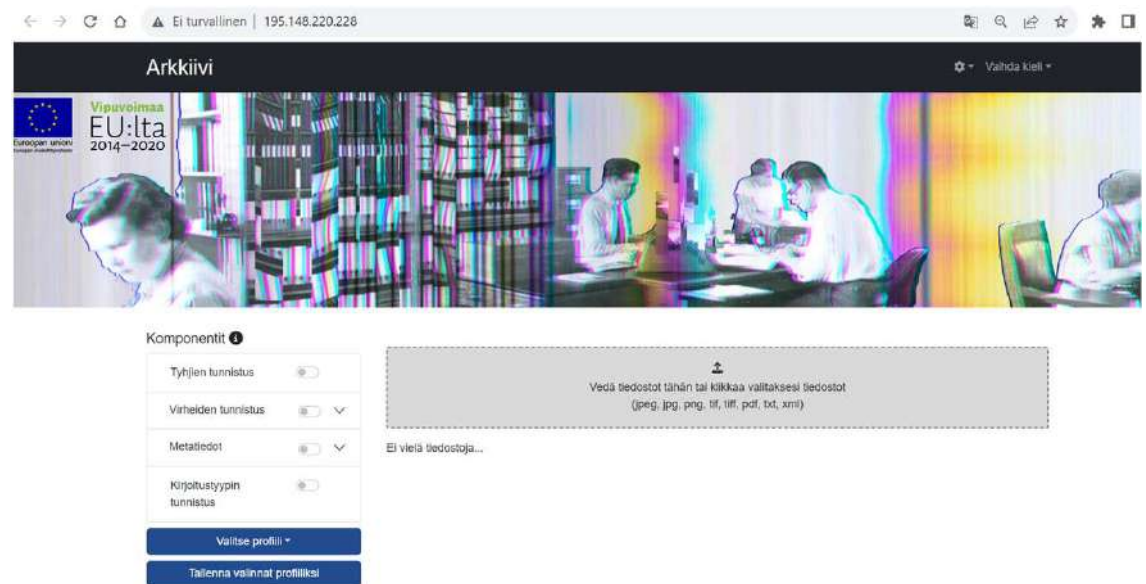
Prebuilt models

Prebuilt models enable you to add intelligent document processing to your apps and flows without having to train and build your own models.



KANSALLISARKISTO

- Kaupallisten pilvipalveluiden haasteena mm. maksullisuus, englannin dominanssi (NLP-sovellukset), pilviprosessoinnin tietosuojahaasteet..
 - Toisaalta isojen toimijoiden resurssit mahdollistavat myös laadukkaiden sovellusten tuottamisen (esim. OCR)
- Myös ilmaisia, suomenkieliselle aineistolle suunnattuja sovelluksia saatavilla lisääntyvässä määrin
- **Arkkiivi** (<https://arkkiivi.fi/>):
 - Alusta sisältää komponentteja jotka tunnistavat digitoiduista aineistoista skannausvirheitä ja sisältöjä
 - Demo/kokeilualusta, joka ei sovellu tuotantokäyttöön: tätä varten sovellusten koodit vapaasti saatavilla GitHubissa





4. Tekoäly Kansallisarkistolla

- 2016-2019 EU-rahoitteinen READ-hanke (Recognition and Enrichment of Archival Documents): Transkribus
- Sisällönanalysointiprojekti 2018-2019 – OCR:aan ja digitoinnin tekoälytyökalujen selvityshanke (DALAI:n esihanke)
- 2020 – Tuomiokirjahaku (3,1 miljoonaa sivua sisältö tunnustettuna)
- 1.9.2021–31.8.2023 DALAI (Digitaalisten aineistojen laadun ja käytettävyyden parantaminen tekoälyavusteisesti)
- 2021 – Pohjoismainen yhteistyö tiivistyy (hackathonit, kuukausittaiset tekoälykehittäjien verkkotapaamiset)
- 2022 -> Oman tekoälykyvykkyyden nosto. Sisällöntunnistuksen ottaminen omiin käsiin ja muiden menetelmien tutkiminen.
- 2023 Uusi Tutkimus ja innovaatiot -toiminto vastaamaan tekoälyn hyödyntämisestä
 - 3 koneoppimisen asiantuntijaa, 5 arkistoalan ja tutkimuksen asiantuntijaa

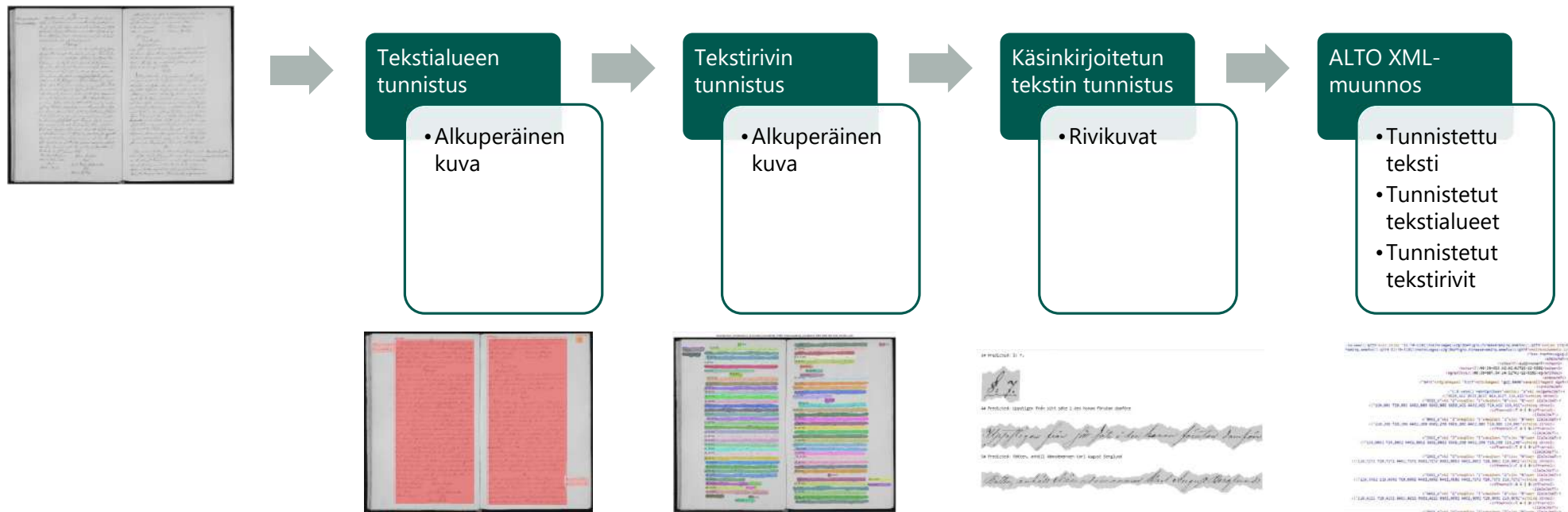
KANSALLISARKISTO

Käsinkirjoitetun tekstin tunnistus

- **HTR** (Handwritten Text Recognition) ollut pitkään haastava osa-alue, käsialojen ja käsinkirjoitetun aineiston monipuolisuuden takia vaikea kehittää yleisiä tekstintunnistusratkaisuja
- Ei juurikaan ilmaisia sovelluksia, vaaditaan usein omien mallien kouluttamista tiettyjä aineistokokonaisuuksia varten
- Arkistojen, kirjastojen ja yliopistojen EU-projektin yhteydessä perustama **Transkribus** (<https://readcoop.eu/transkribus>) on tunnettu ja monien kulttuuriperintöalan toimijoiden käyttämä alusta dokumenttien annotointiin, käsialamallien kouluttamiseen ja hyödyntämiseen
 - Kansallisarkisto ollut myös mukana perustamisesta asti, mutta korkeiden kustannusten takia siirrytty omien HTR-mallien kehittämiseen

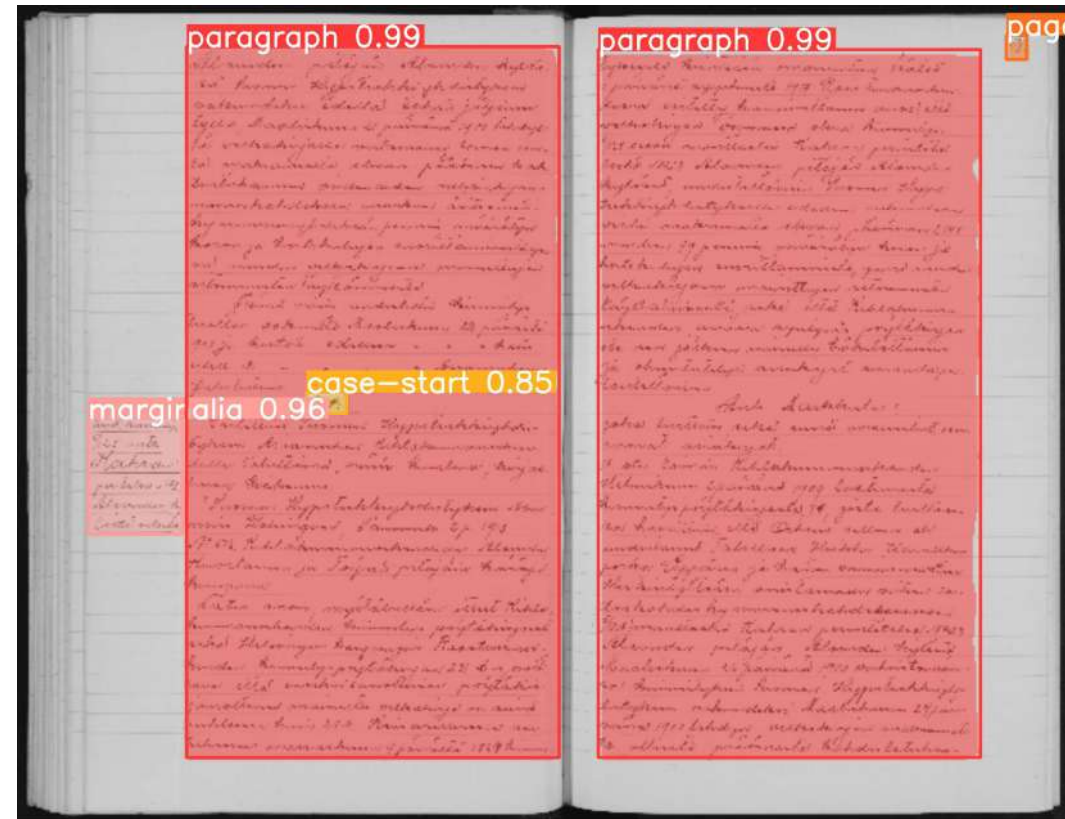
Käsinkirjoitetun tekstin tunnistus

- Esimerkki tekstintunnistuksen eri työvaiheista



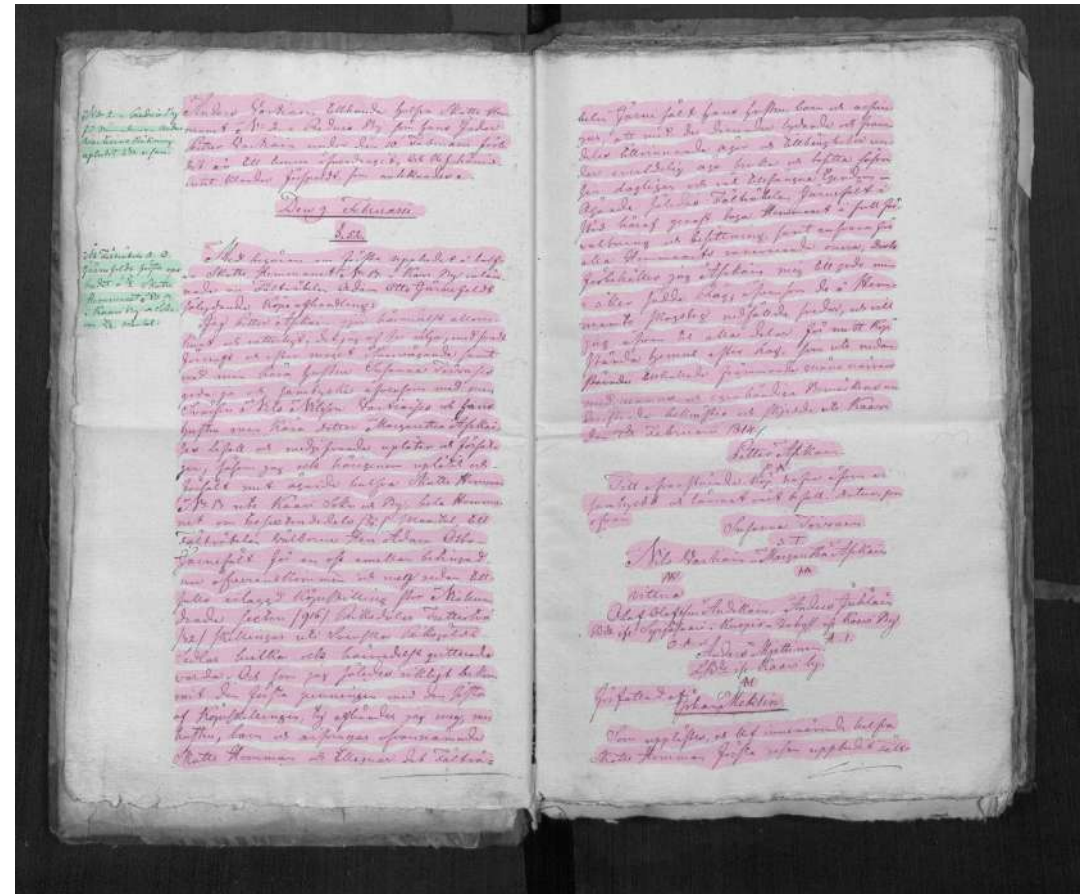
Tekstialueiden tunnistus

- Tunnistetaan sivulta tai aukeamalta eri tyyppisiä kokonaisuuksia
- Tuomiokirjaesimerkissä malli tunnistaa neljän tyyppisiä alueita:
 - Tekstikappale
 - Marginaalissa oleva teksti
 - Sivunumero
 - Uuden oikeusjutun alku



Tekstirivien tunnistus

- Tunnistaa dokumentista tekstirivit
- Rivit yhdistetään koordinaattien perusteella aiemmin tunnistettuihin tekstialueisiin
- Hyödynnetään objektintunnistukseen erikoistunutta neuroverkkomallia



Tekstisisällön tunnistus

- Tekstintunnistusmalli saa syötteenä rivikuvia ja tunnistaa niistä tekstisisältöä
- Neuroverkkomallin koulutusdatana käytetty yli 300 000 rivikuvaa, joihin on käsin annotoitu tekstisisältö

4509# Predicted: van velkomuksen perusteena oleva kont-
Original: van velkomuksen perusteena oleva kont-

van velkomuksen perusteena oleva kont-

4510# Predicted: tokurantti oli oikein uloskirjoitettu Fredriksi
Original: tokurantti oli oikein uloskirjoitettu Fredriks-

tokurantti oli oikein uloskirjoitettu Fredriks

4511# Predicted: sonin reskontrista vuosilta 1887, 1888, 1889 ja
Original: sonin reskontrista vuosilta 1887, 1888, 1889 ja

sonin reskontrista vuosilta 1887, 1888, 1889 ja

header-coller-col header-col header-col

header-coller-col header-col header-col

Opetusaineiston tuottaminen

Lapsi poika virkeä virhosi. Kuoli rasz. aikana, synnyt. aikana, tulla lianditi valekuollut ei vironnut syösti valekuollut.
 synnytyksen jälkeen klo Kuolinsyy
 Pituus 17 sm. Paino 2700 g. Pään ympärys sm.
 Synnytysoireuttomia tai kehityshäiriöitä ei havaittu
 Lapsen vointi myöhemmin (ihottuma, keltatauti y. m.) hyvä

Lisäselvityksiä: (erikoistoinenpiteet, sisätluk., niiden syy, tapauksen kulku, lääkkeit y. m.)
Seikalla vaimalla

- äidin nimi 1
- asuinpaikka 2
- vuosi 3
- täysiaikaisuus 4
- äidin ikä 5
- syntymäpäivämäärä 6
- syntymäkellonaika 7

Sotilaskanta kortti

1. Eklund, Nils Ragnar
 2. Synt. 8. 8. 1917 Helsinki Ukkiniemi 4. Luut.
 5. OKA
 8. Uutonen
 12. Pyhärinta, Rihtniemi
 14. Pyhärinta, Rihtniemi
 16. Pyhärinta, Rihtniemi
 17. Arvostelu: Pyhärinta, Rihtniemi
 19. Palveluksen laus: Pyhärinta, Rihtniemi

Kutunnankilpailu: 247
 a. Kutunnankilpailu: 247
 b. Vapantettu: 19
 i. Palvelusaika: 1917-1941
 j. Vakkainan aavinpaikka ja osoite kotintam- sen jälkeen: Nil

Suolaks by, af Kronolänsmannen Anders Hagman, å, tjenstens vägnar samt Torparen Henrik Passila, från sagde by, i egen skap af målsegande under tilltal ställd för delaktighet ute

- 1-1 384
- 2-1 Suolaks by, af Kronolänsmannen Anders Hagman, å, tjenstens
- 2-2 vägnar samt Torparen Henrik Passila, från sagde by, i egen
- 2-3 skap af målsegande under tilltal ställd för delaktighet ute

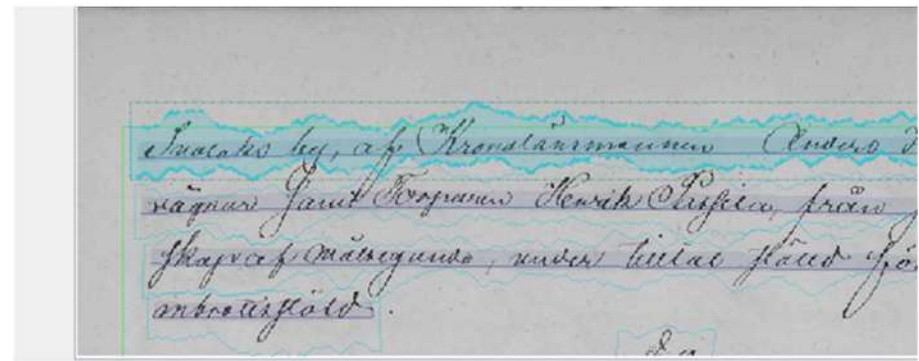
1 ~~lijffztijdh Eij kunnat skaffa H:r Burkhardtz Consens, till kiöpet~~
 2 ~~af Erich Afwin för 3:ne tunnor Tiära förledhen Sommars fördes~~
 3 ~~på Klimpos~~
 4 ~~Lådia, till upfylla dhen andre Påst Tiära Erich Afwinen till hörigh~~
 5 ~~war~~
 6 ~~ophäfwes, effter andre borgare seijia sikh aldrigh fracht begiära~~
 7 ~~häller taga för~~
 8 ~~fyllningh Tiära, hälst som han Eliest så dyrr fracht och utföore för~~
 9 ~~tunna~~
 10 ~~tagit.~~
 11 ~~Resolverades att Matz Juoin betahlar för sit ifrån Stockholm~~
 12 ~~kommandhe~~

7 Enckia betahla tillbakas äth Pafwillla 3 D:r, effter Sahl:gh
 8 Torfwinen i sin
 9 lijffztijdh Eij kunnat skaffa H:r Burkhardtz Consens, till kiöpet
 10 som han
 11 läfwadhe, men 2 D:r behåller Enckian för dhen mödan hannes
 12 Sahl:ge Man
 13 brukat till förmå mehr bem:e H:r Burkhardtz tillståndh att
 14 städfästa kiöpet,
 15 derföre Han och måst någre reesor Påstpeningar betahla.
 16 Huadh wedh kommer om dhe 2 D:r 8 / öre som Staphan Klimpo
 17 fordrar i fracht
 18 af Erich Afwin för 3:ne tunnor Tiära förledhen Sommars fördes
 19 på Klimpos
 20 Lådia, till upfylla dhen andre Påst Tiära Erich Afwinen till hörigh
 21 war
 22 ophäfwes, effter andre borgare seijia sikh aldrigh fracht begiära
 23 häller taga för
 24 fyllningh Tiära, hälst som han Eliest så dyrr fracht och utföore för
 25 tunna
 26 tagit.
 27 Resolverades att Matz Juoin betahlar för sit ifrån Stockholm
 28 kommandhe

margin
 1
 2

Kansallisarkiston opetusaineistokokonaisuuksia

- READ-hanke & 1800-luvun tuomiokirjat: n. 1,3 miljoonaa sanaa
- Joukkoistushankkeet: n. 2 miljoonaa sanaa
 - Sotapäiväkirjoja
 - 1800-luvun perukirjoja
 - Henkikirjoja
- Muut arkistot ja tutkijat: n. 3 miljoonaa sanaa
 - Riksarkivetilta 1600- ja 1800-lukujen aineistoa
 - Ville Pekka Kääriäinen (HY) merkittävä määrä 1600-lukua
 - Tukholman kaupunginarkistolta 1700-luvun aineistoja

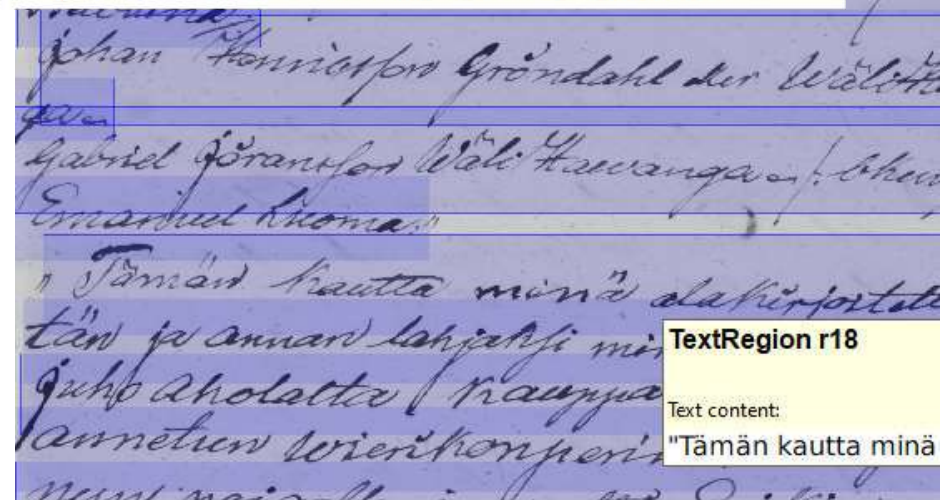


1-1 384

2-1 Suolaks by, af Kronolänsmannen Anders Hagman, å, tjenstens

2-2 vägnar samt Torparen Henrik Passila, från sagde by, i egen

2-3 skap af målsegande, under tilltal ställd för delaktighet ute



TextRegion r18

Text content:

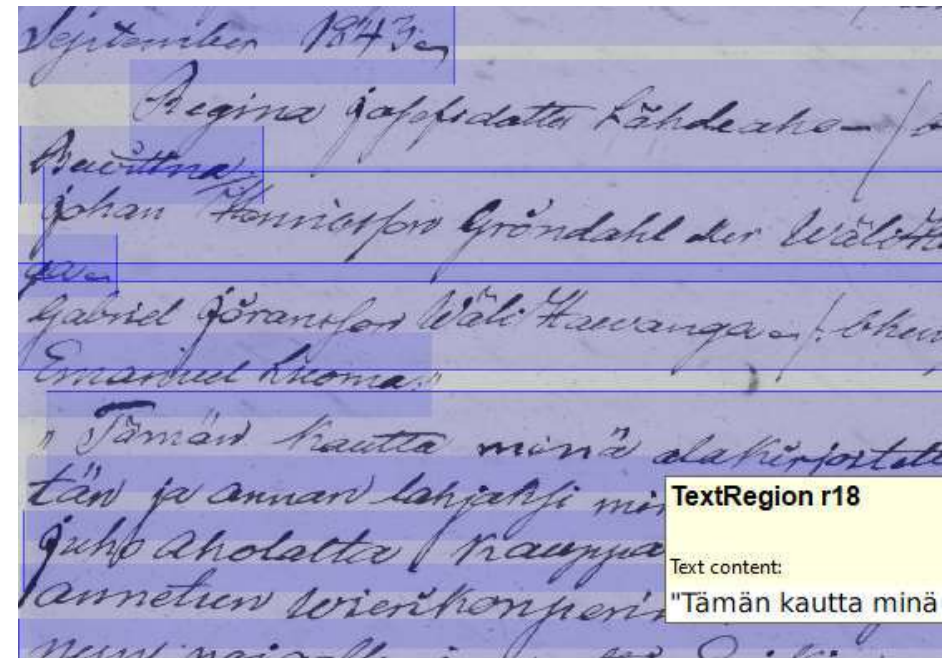
"Tämän kautta minä

KANSALLISARKISTO

Tuomiokirjat-malli

- Tuomiokirjojen käsialamalli perustuu kahteen tekoälymalliin
- Segmentointi: YOLOv8
 - Parannettu segmentointia, mm. juttujen alkujen tunnistus
- Tekstintunnistus: TrOCR
 - Microsoftin kehittämä käsiala-aineistolla esikoulutettu malli
 - Opetusaineistoa yli miljoona sanaa

Merkkivirhe 2,4 %



4509# Predicted: van velkomuksen perusteena oleva kont-
Original: van velkomuksen perusteena oleva kont-

van velkomuksen perusteena oleva kont-

4510# Predicted: tokurantti oli oikein uloskirjoitettu Fredriksi
Original: tokurantti oli oikein uloskirjoitettu Fredriks-

tokurantti oli oikein uloskirjoitettu Fredriks-

4511# Predicted: sonin reskontrista vuosilta 1887, 1888, 1889 ja
Original: sonin reskontrista vuosilta 1887, 1888, 1889 ja

sonin reskontrista vuosilta 1887, 1888, 1889 ja

KANSALLISARKISTO

Suomalainen supermalli v. 0.1

- Viimeisin kokeilu käsialantunnistukseen liittyen
- Käsialamalli, joka lukee tekstiä 1600-luvulta 1920-luvulle
 - Koulutettu yhteensä 737 000 rivillä
 - Laskenta kesti reilun kaksi viikkoa
 - Virheprosentti 3,66
- Opetusaineistoa kerättynä huomattava määrä lisää
 - Supermallin seuraava versio vuoden 2024 aikana
 - Huomattava määrä lisäaineistoa -> tarkkuuden parantuminen
 - Jatkokoulutus 1500-luvulle?
- Supermalli toimii myös pohjamallina
 - Koulutettu useita aineistokohtaisia käsialamalleja supermallin päälle -> ei tarvita kovin paljoa uutta aineistoa

191# Predicted: at tå han, uppå århållen kal-
Original: at tå han, uppå århållen kal-

at tå han, uppå århållen kal-

192# Predicted: telse af Inspektoren Livijn, wa-
Original: telse af Inspektoren Livijn, wa-

telse af Inspektoren Livijn, wa-

193# Predicted: rit til Spinhuset at underrätta
Original: rit til Spinhuset at underrätta

rit til Spinhuset at underrätta

44321# Predicted: bestridt alla honom åliggande Presterliga
Original: bestridt alla honom åliggande Presterliga

bestridt alla honom åliggande Presterliga

44322# Predicted: görenal till Consistorii och Församlingens
Original: göremål till Consistorii och Församlingens

göremål till Consistorii och Församlingens

44323# Predicted: ofullkomliga nöje, under iakttagande
Original: fullkomliga nöje, under iakttagande

ofullkomliga nöje, under iakttagande

4509# Predicted: van velkomuksen perusteena oleva kont-
Original: van velkomuksen perusteena oleva kont-

van velkomuksen perusteena oleva kont-

4510# Predicted: tokuranti oli oikein uloskirjoitettu Fredriksi
Original: tokuranti oli oikein uloskirjoitettu Fredriks-

tokuranti oli oikein uloskirjoitettu Fredriks-

4511# Predicted: sonin reskontrista vuosilta 1887, 1888, 1889 ja
Original: sonin reskontrista vuosilta 1887, 1888, 1889 ja

sonin reskontrista vuosilta 1887, 1888, 1889 ja

KANSALLISARKISTO

191# Predicted: at tå han, uppå ärhållen kal-
Original: at tå han, uppå ärhållen kal-

A horizontal strip of aged, yellowish paper with handwritten text in a cursive script. The text is written in dark ink and appears to be a Swedish phrase, matching the predicted and original text above it.

192# Predicted: telse af Inspektoren Livijn, wa-
Original: lelse af Inspektoren Livijn, wa-

A horizontal strip of aged, yellowish paper with handwritten text in a cursive script. The text is written in dark ink and appears to be a Swedish phrase, matching the predicted and original text above it.

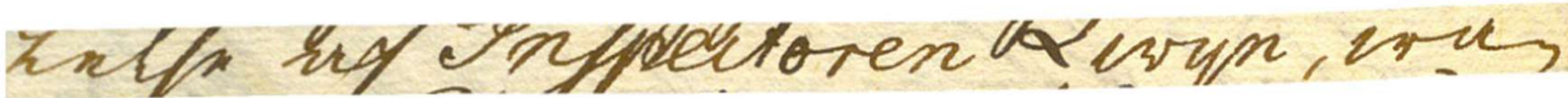
193# Predicted: rit til Spinhuset at underrätta
Original: rit til Spinhuset at underrätta

A horizontal strip of aged, yellowish paper with handwritten text in a cursive script. The text is written in dark ink and appears to be a Swedish phrase, matching the predicted and original text above it.

191# Predicted: at tå han, uppå ärhållen kal-
Original: at tå han, uppå ärhållen kal-

A horizontal strip of aged, yellowish paper with handwritten text in dark ink. The text is written in a cursive script and appears to be a Swedish phrase: "at tå han, uppå ärhållen kal-".

192# Predicted: telse af Inspektoren Livijn, wa-
Original: lelse af Inspektoren Livijn, wa-

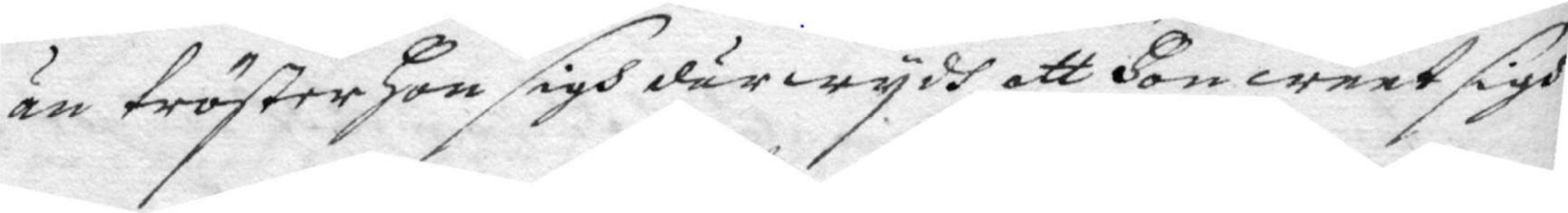
A horizontal strip of aged, yellowish paper with handwritten text in dark ink. The text is written in a cursive script and appears to be a Swedish phrase: "telse af Inspektoren Livijn, wa-".

193# Predicted: rit til Spinhuset at underrätta
Original: rit til Spinhuset at underrätta

A horizontal strip of aged, yellowish paper with handwritten text in dark ink. The text is written in a cursive script and appears to be a Swedish phrase: "rit til Spinhuset at underrätta".

83# Predicted: än tröster hon sigh där wydh att hon weet sigh

Original: än tröster Hon sigh därwydh att Hon weet sigh



än tröster hon sigh därwydh att hon weet sigh

84# Predicted: wara oskyldigh, och tackar Gudh därföre

Original: wara oskyldigh, och tackar Gudh därföre



wara oskyldigh, och tackar Gudh därföre

1500-luku!

A horizontal strip of a handwritten manuscript in a cursive script, likely from the 1500s. The text is partially visible and appears to be in a historical language.

GT: Riddare, att han medh edher trogne undersåtherr, så
Pred: Riddare, att han medh edher wogne undersåther, så
CER : 0.059

Another horizontal strip of a handwritten manuscript, similar to the first one, showing cursive handwriting on aged paper.

GT: tree eller fyre präster, effter som the haffue gelden och
Pred: rer eller fyra präster, effter som dhe haffue gatan och
CER : 0.14

Rakenteelliset aineistot

- Kansallisarkisto täynnä erilaista vanhaa tilasto- ja rekisteriaineistoa
 - Henkikirjat ja kirkonkirjat
 - Väestönlaskentatietoja
 - Kuolinsyytilastot
 - Laivojen lokikirjoja
 - Kiinteistörekisterit
- Aineistojen koneellinen tunnistus tarjoaa huikeat mahdollisuudet tuottaa aineistoja laajasti eri tieteenoaloille
- Taulukkorakenteiden tunnistus ollut pitkään ongelma tekoälylle, mutta keinoja alkaa nyt löytyä

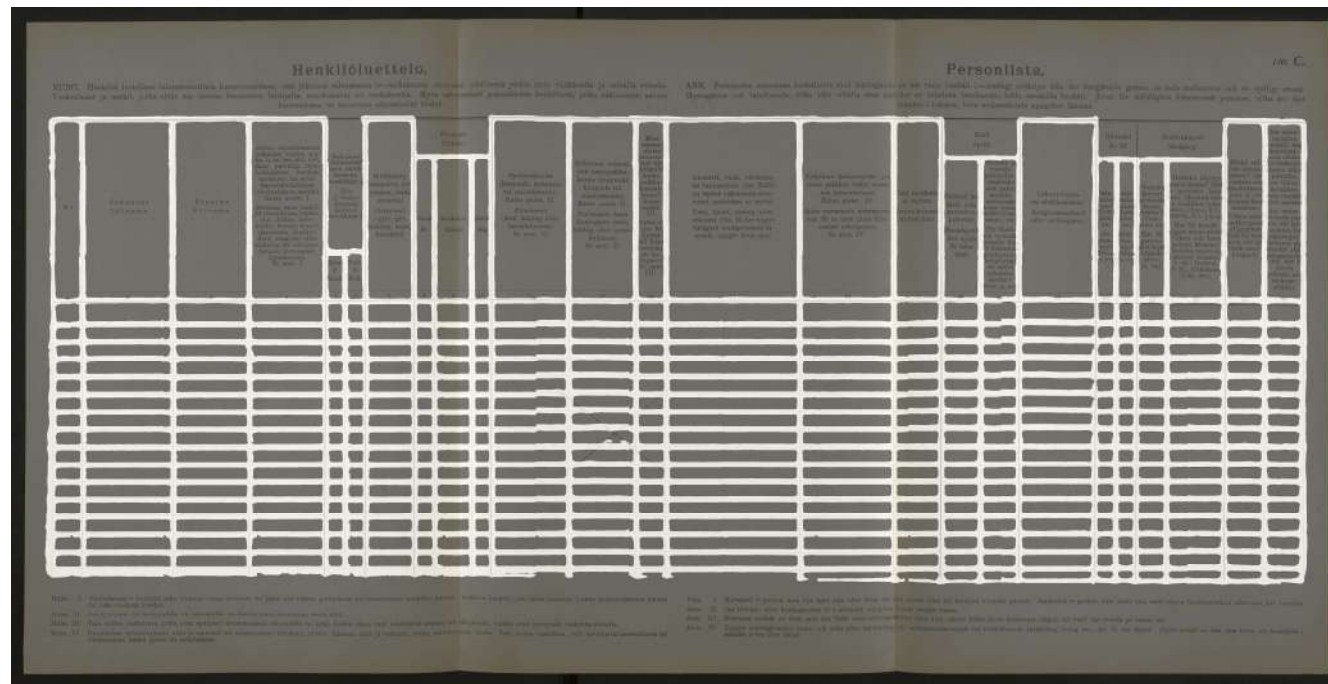
Page 9 Haiholais Natthöhl.

Haiholais by.	Födelse		Komman itrin	Köpen	Läsning utur minnet					Födelse Ort	Vid läsning stället	Natthvardsgång				
	År och Datum	Ort			År	Ort	År	Ort	År			Ort	18/6	18/7	18/8	18/9
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						
...	X	X	X	X	X	X						

KANSALLISARKISTO

Taulukkomaskin tunnistus

- Annotoiduilla taulukkokuvilla koulutettu neuroverkkomalli tunnistaa kuvasta taulukkorakenteen
 - Vaatii 200-300 annotoitua taulukkoa toimiakseen
- Lopputuloksena taulukkosolujen koordinaatit



Tekstin tunnistus taulukkosoluista

- Taulukkokoordinaattien perusteella solut voidaan leikata omiksi kuviksi
- Solukuvat toimivat syötteenä tekstintunnistukselle
- Tuloksena ennustettu tekstisisältö, joka voidaan viedä esim. excel-taulukkoon



KANSALLISARKISTO

50025 17/4-36

TILASTOLLINEN PÄÄTOIMISTO

16.4.2024

Kuolleet Haminan seurakunnassa v:n 1936 I neljänneksellä.

(Lähetettävä Tilastolliseen päätoimistoon 15 päivän kuluessa neljännessuoden päättymisestä)

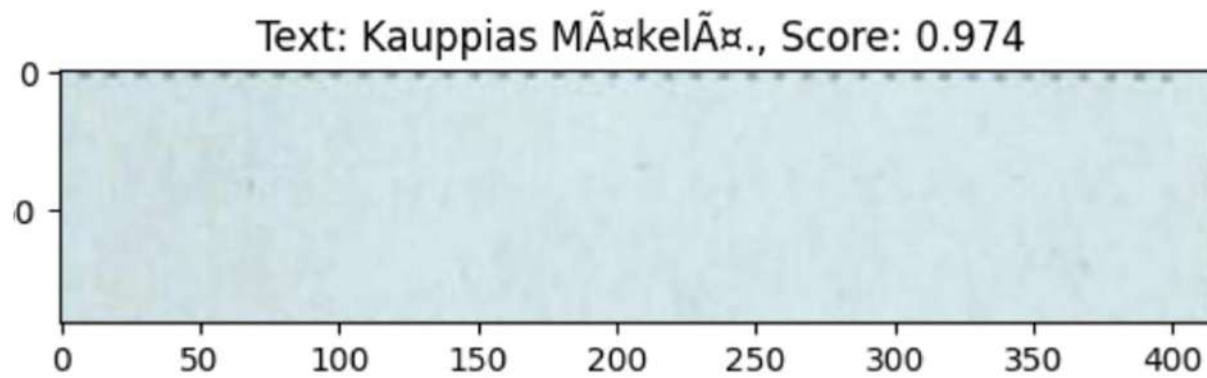
Kuoleman aika	Sivu pää-	Kuolleen nimi, sääty ja osoite	Syntymä- aika	Ikä			Naimen, Leski, Er- Lapsi	Kuoleman syy	Miehen	Vaino
				V	K	P				

	A	B	C	D	E	F	G	H	I	J	K	L
1	Kuukausi	Päivä	u pääkirja	Kuolleen nimi, sääty ja osoite	Syntymäaika	Ikä: Vuotta	Kuukaut	Ikä: Päivää	ut, Leski,	Kuoleman syy	Miehenp.	Vaimonp.
2	Tammik.	12	591	Saikkonen, Iida Maria, työm. vo	7.I 1884	52		5	Nainut	a) Carcinoma cordci e anemia secun	dis	1
3	"	15	546	Rautaleins, Matilda, johtajan leski	13.III. 1876	59	10	2	leski	a) Myodegenera tio cordis (3030 sar	81	lis
4	"	16	920	Lautala, Antti, kalastaja	22.XII	95		24	"	a) Haama- raspia cordis	1	
5	"	20	371	Lindblad, Elina Augusta, työnjoht.	2.IX. 1859	76	4	14	"	a) Arterio sc b) Marasmus sc	eron	1
6	"	27	428	Nakari, Tuomas, suutari	1.II 1873	62	10	26	n-nut	Harmorrhag cerebri. b) Samoin.	19/	
7	"	29	297	Korttila, Hellin, laivankulj. tytär.	27.VII 1914	21	6	2	n:ton	a) Tuberculm pulmo. b) sa	s. mini	1
8	Helmik.	5	512	Pousi, Anna Maria, leski	20.I. 1867	68		15	leski	a) Myodegener et insufficien cordis atio tio siv	1	
9	"	13	689	Takalahti, Kaarle Antinp. talonmies.	25.I. 1865	71		18	nton	a) Myodegene tio et insuffia cordia sa 7/5 ole in.		
10	"	17	44	Betán, Mariaana, leski	20.V 1874	61	8	27	leski	a) Vitium ca 3020; b) Insu cientia cor	Pelin	1
11	"	13	71	Eriksson, Aron Johan, ent. pol. konst	15.VI 1876	59	7	28	n-nut.	a) Hypertonio arterialis. mon	1	
12	"	25	290	Kitunen Helena. leski	6.1. 1849	87	1	20	leski	a) Marasmus seonlis; b.) Bro pneum	velko-	1
13	Maalisk	19	426	Naakka, Amanda, työmiehen vaimo	6.9.1898	37	6	13	n-nut	a) Tuberculo meningrum	sia. bää in	1
14	"	23	73	Eriksson, Vilhelm. työmies	6.8 1862	73	7	17	"	a) Samoo kepatio. b) samoin.	1	
15	"	28	464	Paavola, Maria Kristiina, leski	23.5. 1851	84	10	5	leski	a&b) Marasm senilis. Myio insuf	-"- lago	lin

"	25	290	Kitunen, Helena, leski	6.1. 1849	87	1	20	leski	a) Marasmus seonlis; b.) Bro pneum meningrum.	1	
Maalisk	19	426	Naakka, Amanda, työmiehen vaimo	6.9. 1898	37	6	13	n-nut	a) Tuberculo meningrum et insu cientia cordis b) samoin.	1	
"	23	73	Eriksson, Vilhelm, työmies	6.8 1862	73	7	17	"	a) Samoo kepatio. b) samoin.	1	
"	28	464	Paavola, Maria Kristiina, leski	23.5. 1851	84	10	5	leski	a) Marasmus senilis. Myio insuf.	1	

LISARKISTO

6 9 Käännä!



Tietopyyntöjen tehostaminen

- Tuhansia tietopyyntöjä ajoneuvorekisterin aineistoihin liittyen vuosittain -> vastaaminen vaati paljon työaika
- Sisällöntunnistuksella apua tietopyyntötyöhön
 - Kuvia annotoitiin n. 800
 - Yolo segmentointimallina, TrOCR tekstintunnistukseen
- **Tehostaa työtä huomattavasti, arviolta laskee tietopyyntöön vastaamisen kuluvan työajan alle puoleen entisestä**

rekisterinumero: R-31
VCE-439 LAJI: 1B, PAKETTI-AUTO MERKKI: FORD VALM.NO: GB71RC11417
***** REKISTERÖITY: 09.04.75 KÄYTTÖNOTTOVUOSI: 75 VAKUUTUS: SVENSK-FINLÄN
MALLI: TRANSIT 100-71E2X-2.4L DIESEL VAN/2690 KÄYTTÖ: YKSITYINEN/1 V.
RENKAAT: 7.50-14/6PR 195R14 AKS.LKM: 2 JÄRJEL: NAP MUUTOSKATS:
- OHJASTAJA: ANDELSLÄGET ÖSTERBOTTENS KÖTT 65100 VASA 10, SMEDSBYVÄGEN 7
KÄYTTÖVOIMA: DIESEL/3 PAINOT: 1470 +980 =2450 KG HENKILÖLUVUT: KL=2 ,KV:
KOTIPAikka: 70/4401/905-0 TYYPPIID: 473043-0 KORITAKENNE: UMPI PELTIKORI
AKSELIPAINOT: 1270 KG KORIYLITYS: 100/ CM LEVEYS: 196/ CM
--AUTOVEROVAPAA PAKETTI-AUTONA, LAKI 352/71 /

rekisterinumero: R-31
VCE-765 LAJI: 1A, HENKILÖAUTO MERKKI: OPEL VALM.NO: B12648801
***** REKISTERÖITY: 03.04.75 KÄYTTÖNOTTOVUOSI: 71 VAKUUTUS: SVENSK-FINLÄN
MALLI: 2D ASCONA 165-A-81/2430 KÄYTTÖ: YKSITYINEN/1 V.
Ajoneuvorekisterin etsintätyökalu
Hakusanat:
311878
Tulokset:
Normaali haku
10
Etsi
Etsitään normaalilla haulla sanoilla: 311878
Dokumentin nimi:
2
Tiedostonimi:
VCE439 0355
OCR:n tulokset kyseisessä dokumentissa:
TS-370 311878 TS-370 311878
Dokumentin nimi:
1
Tiedostonimi:
298639437_0956
OCR:n tulokset kyseisessä dokumentissa:
311878 TS-370
Etsimiseen käytetty aika: 2.62 sekunttia

KANSALLISARKISTO



226 m

Digitaalista kuvaa



~62 m

OCR



3,5 m

HTR

KANSALLISARKISTO

Aineistojen prosessointisuunnitelmia

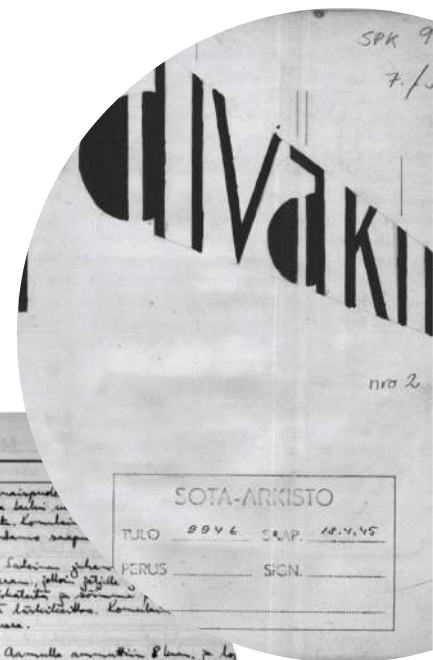
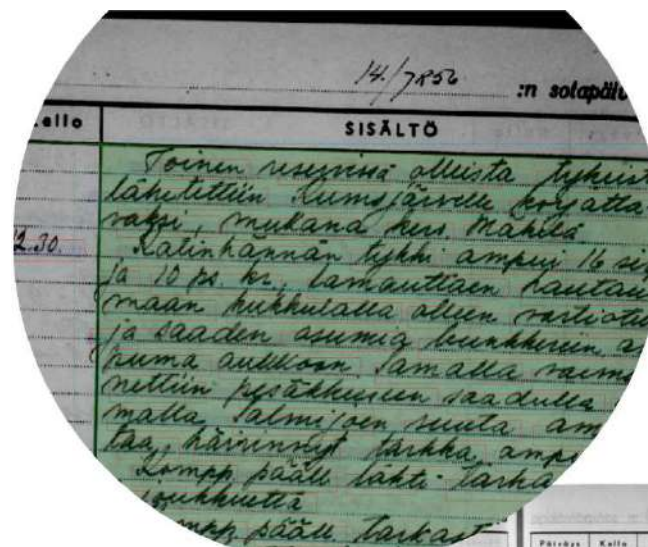
- KA päässyt mukaan CSC:n akateemisen käytön piiriin
 - Merkittävät laskentaresurssit käyttöön
 - Jatkossa jopa 20 miljoonaa sisältötunnistettua sivua vuosittain
- Parhaillaan prosessoidaan 1800-luvun tuomiokirja-aineistoa
 - Yhteensä hieman alle 10 miljoonaa sivua
 - Käytännössä koko 1800-luvun oikeushistoria sisältöhakujen piiriin vielä tämän vuoden aikana!
- 1800-luvun valmistuttua siirryttään 1700-luvulle ja ensin renovoituihin tuomiokirjoihin



KANSALLISARKISTO

Sotapäiväkirjat

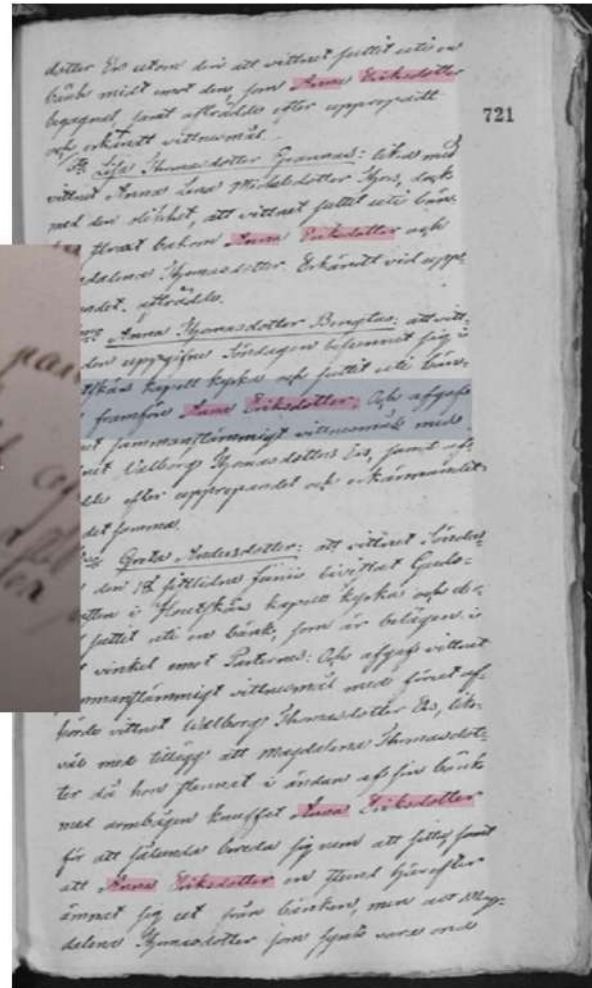
- Kansallisarkiston suurin joukkoistushanke, alkoi keväällä 2022
- Yli 200 vapaaehtoista, joista yli 100 osallistui aktiivisemmin
- Transkribus Lite –työkalun avulla sotapäiväkirjojen puhtaaksikirjoitusta
- Yli miljoona sanaa opetusaineistoa
 - Vaatinut kuitenkin paljon korjailua ja läpikäymistä
- Nyt kuitenkin valmista, ja prosessointi aloitetaan kesään mennessä
 - Taulukkomalli, joka jakaa tapahtumat päivän mukaan



KANSALLISARKISTO

Search Finnish Court Records
 Search and browse Finnish court records with the help of algorithmic text recognition.

Search



dotter Er utom den att vittnet suttit uti en bank midt emot den, som Anna Eriksdotter begagnat samt afträdde efter upprepadt och erkändt vittnesmål.

6o Lisa Thomasdotter Grannas: likd med vittnet Anna Lena Michelsdotter Thors dock med den denhet att vittnet suttit uti banken straxt bakom Anna Eriksdotter och Magdalena Thomasdotter. Erkändt vid upprepadet. afträdde.

1mo Anna Thomasdotter Bengtas: att vittnet den uppgifne Föndagen befunnit sig i Houtskän kapell kyrka och suttit uti banken framföre Anna Eriksdotter; Och afgaf vittnet sammanstämmigt vittnesmål med vittnet Valborg Thomasdotters Ers, samt afträdde efter upprepadet och erkännandet af det samma.

8° Greta Andersdotter: att vittnet Söndagen den 18. sistlidne Junii bivistat Gudotjensten i Houtskäns kapell kyrka och dervid suttet uti en bank, som är belägen i rat vinkel emot Parternes; Och afgaf vittnet sammanstämmigt vittnesmål med förut afporde vittnet Walborg Thomasdotter Ers, liksom

KANSALLISARKISTO

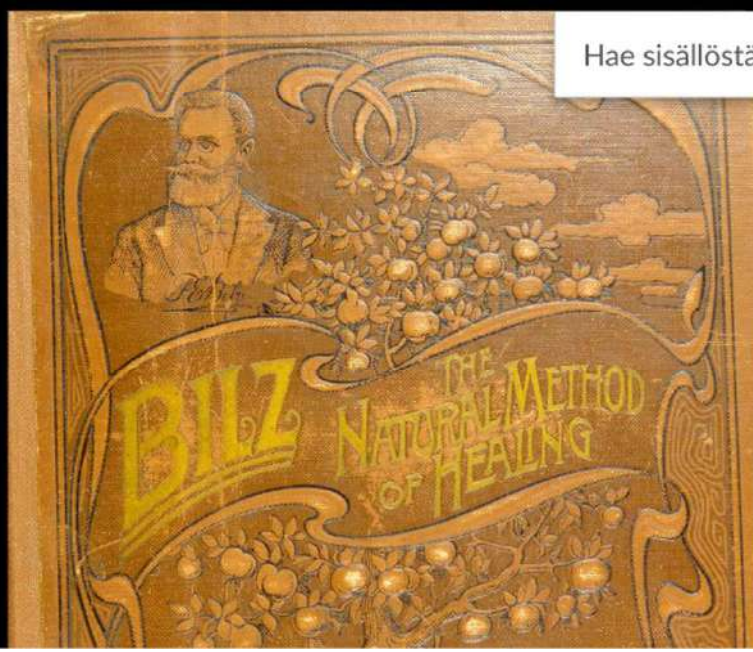
The natural method of healing: A new and complete guide to health by F.E. Bilz

Ruudukko

Viitetiedot

Sulje haku x

>>



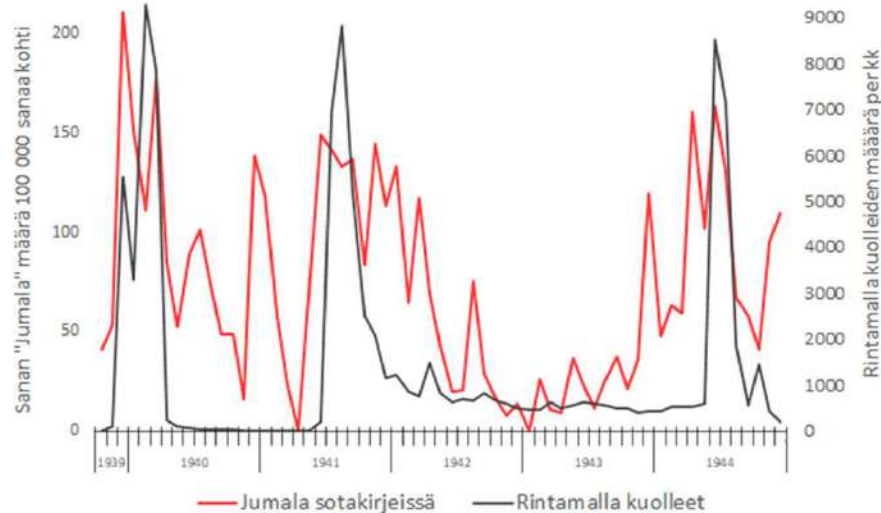
Hae sisällöstä tekstitunnistuksen avulla

Hae Q

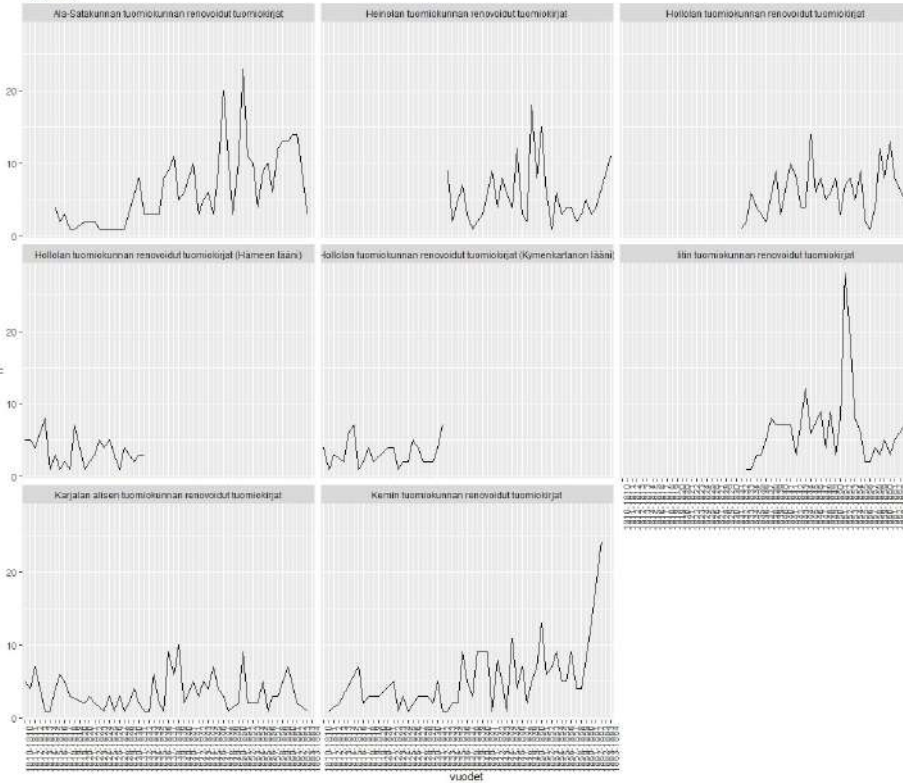


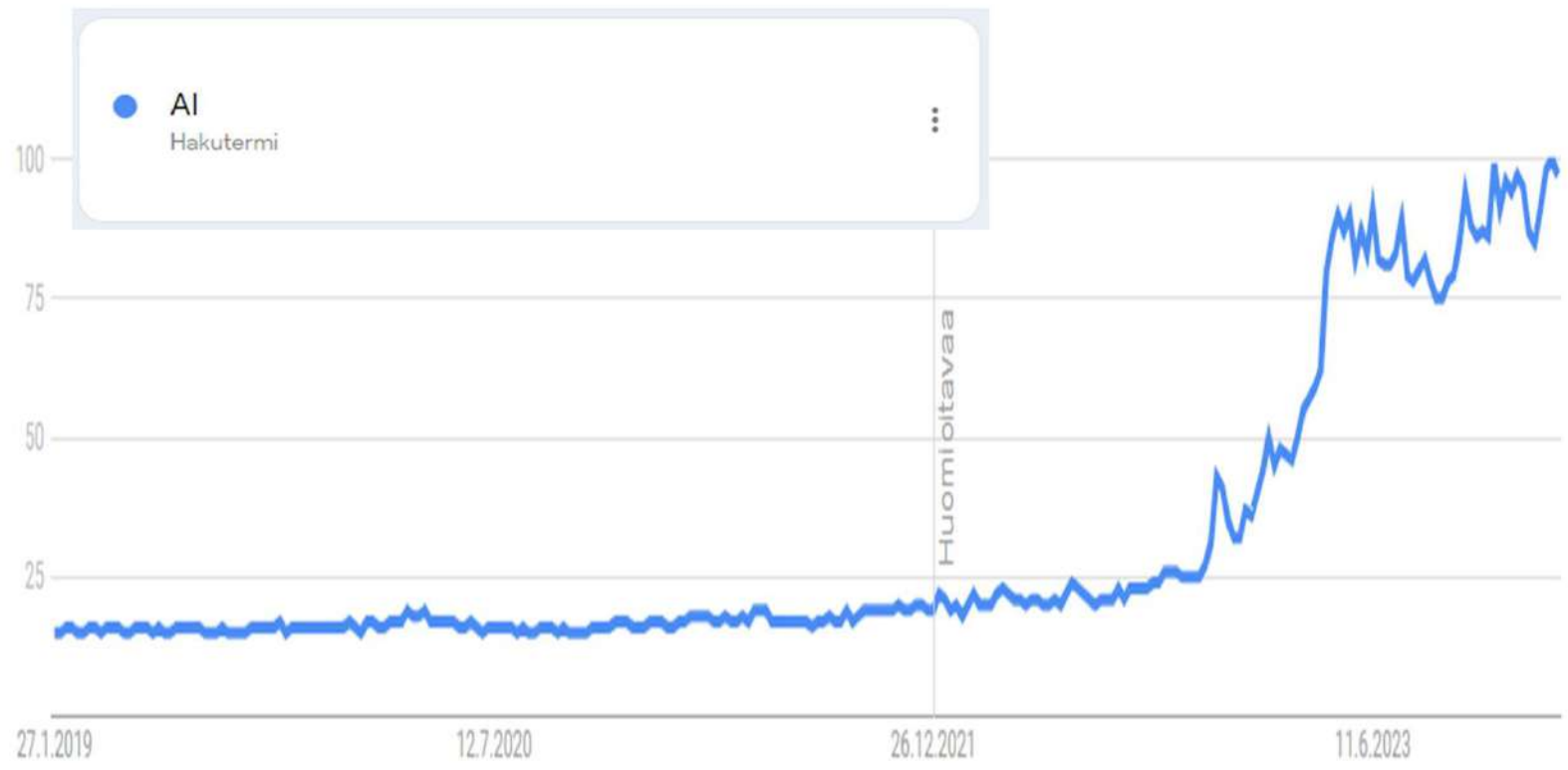
Merkitys tutkimukselle

- Koneluettavuus mahdollistaa uudet menetelmät
- Taulukkomuotoisen aineistot pian mahdollista tuottaa suoraan laskennalliseen tutkimukseen kelpaavaksi dataksi



Kasvatillaisiin ja elatukseen liittyvät oikeusjutut
1810-1864





Kielimallien mahdollisuudet

- Suurten kielimallien potentiaalia asiakirjanhallintaan liittyvissä sovelluksissa ollaan aktiivisesti tutkimassa eri tahoilla
- Mahdollisia käyttökohteita: tiivistelmien teko tekstikokonaisuuksista, kielimallien yhdistäminen asiakirja-aineistoon niin että ne vastaavat käyttäjän aineistoa koskeviin kysymyksiin, kielimallien hyödyntäminen aineistohaussa (esim. vektoritietokannat)...
- Kielimallien käyttöön liittyy kuitenkin monia haasteita:
 - Luotettavuus
 - Maksullisuus
 - Kielivalikoima
 - 'Raakamallit' vs. eri käyttötarkoituksiin jalostetut mallit
 - Resurssivaatimukset...
- Parhaillaan kuitenkin kehitteillä avoimia, suomenkielisellä aineistolla koulutettuja malleja
 - TurkuNLP:n GPT-mallit, TukuNLP:n ja Silo AI:n Poro
 - Näiden myötä toivotaan ratkaisuja osaan haasteista



[TurkuNLP/gpt3-finnish-large](#)



KANSALLISARKISTO

```

... index=index)
bot.load_chat_history("chat_history.json")

while True:
    user_input = input("You: ")
    if user_input.lower() in ["bye", "goodbye"]:
        print("Bot: Goodbye!")
        bot.save_chat_history("chat_history.json")
        break
    response = bot.generate_response(user_input)
    print(f"Bot: {response['content']}")

```

... You: Milloin suomessa oli sisällissota ja mistä löydän siitä tietoa

Bot:

Suomessa oli sisällissota vuonna 1918. Tietoa sisällissodasta voi löytää Sisällissodan pikkujättiläinen -teoksesta, joka on julkaistu vuonna 2009. Lisäksi Suome

You: Mitä arkistoaaineistoja sisällissodasta on?

Bot:

Vapaussodan arkisto sisältää hyvin laajasti aineistoa kansalaissodasta, mukaan lukien punaisen ja valkoisen puolen organisaatioiden tuottamia asiakirjoja, yksit

You: Entä muita?

Bot:

Muita lähteitä, joista voi etsiä tietoa valkoisen armeijan puolella toimineista henkilöistä, ovat kunniamerkkivaliokunnan arkisto ja erilaiset valkoisen puolen i

You: Mistä löydän tietoa punaisista?

Bot:

Vapaussodan arkistosta voi löytää hyvin monenlaista tietoa niin sotilaallisista kuin siviilipuolen asioista. Aineisto on jaettu arkistossa keräyspiirin mukaan j

Vapaussodan arkistoa voi tutkia Kansallisarkiston Helsingin toimipisteessä. Sieltä löytyvät myö

You: Mistä löydän tietoa sisällissodassa taistelleista punaisista?

Bot:

Tietoa sisällissodassa taistelleista punaisista voi etsiä Suomen Vapaussota vuonna 1918 -teoksen ja Suomen Vapaussota -teoksen lisäksi muista lähteistä, kuten s

You: Mistä voin löytää tietoa yksittäisestä punaisten puolella taistelleesta henkilöstä?

Bot:

Tietoa yksittäisestä punaisten puolella taistelleesta henkilöstä voi etsiä punaisten puolella taistelleiden joukko-osastojen arkistoista. Kansallisarkistossa He

You:

KANSALLISARKISTO

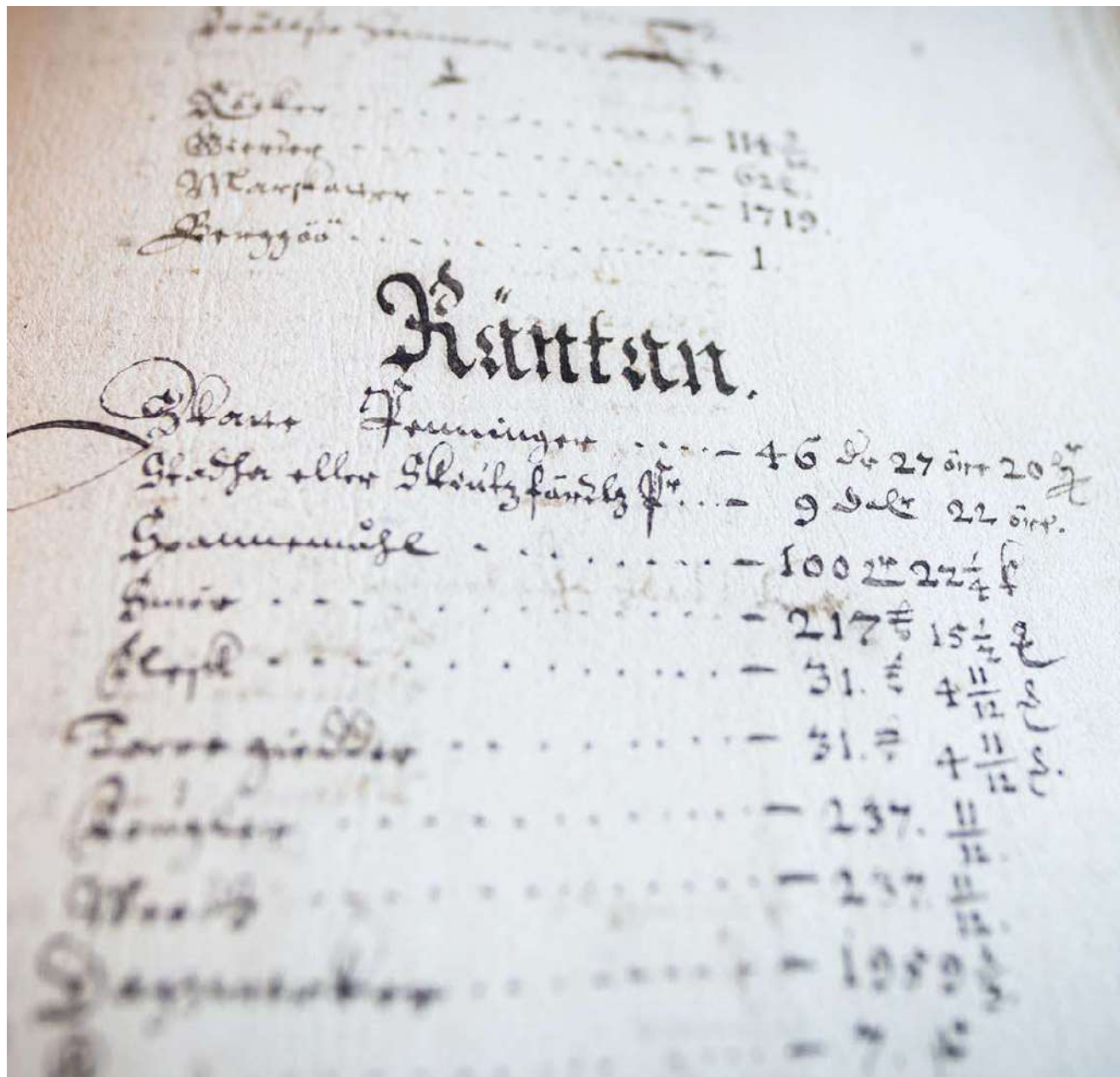


KANSALLISARKISTO

Kiitos!

Ilkka.jokipii@kansallisarkisto.fi

mikko.lipsanen@kansallisarkisto.fi



Esitys on toteutettu osana Kaakkois-Suomen ammattikorkeakoulu Xamk:n, Ammattiopisto Samiedun ja Kansallisarkiston yhteishanketta JoDi-Joustavat koulutus ja työelämäpolut tulevaisuuden digitointiosaajille (ESR+ 2023-2025).



**Euroopan unionin
osarahoittama**





KANSALLISARKISTO

www.kansallisarkisto.fi



@kansallisarkisto



@kansallisarkist