# Practical Open Source AI "stubs" for archives – PoC solutions

*Ph.D Anssi Jääskeläinen*

*Xamk / Digitalia reseach center*

*Mikkeli, Finland*

# Agenda: From files to archive with AI
## Set of AI tools

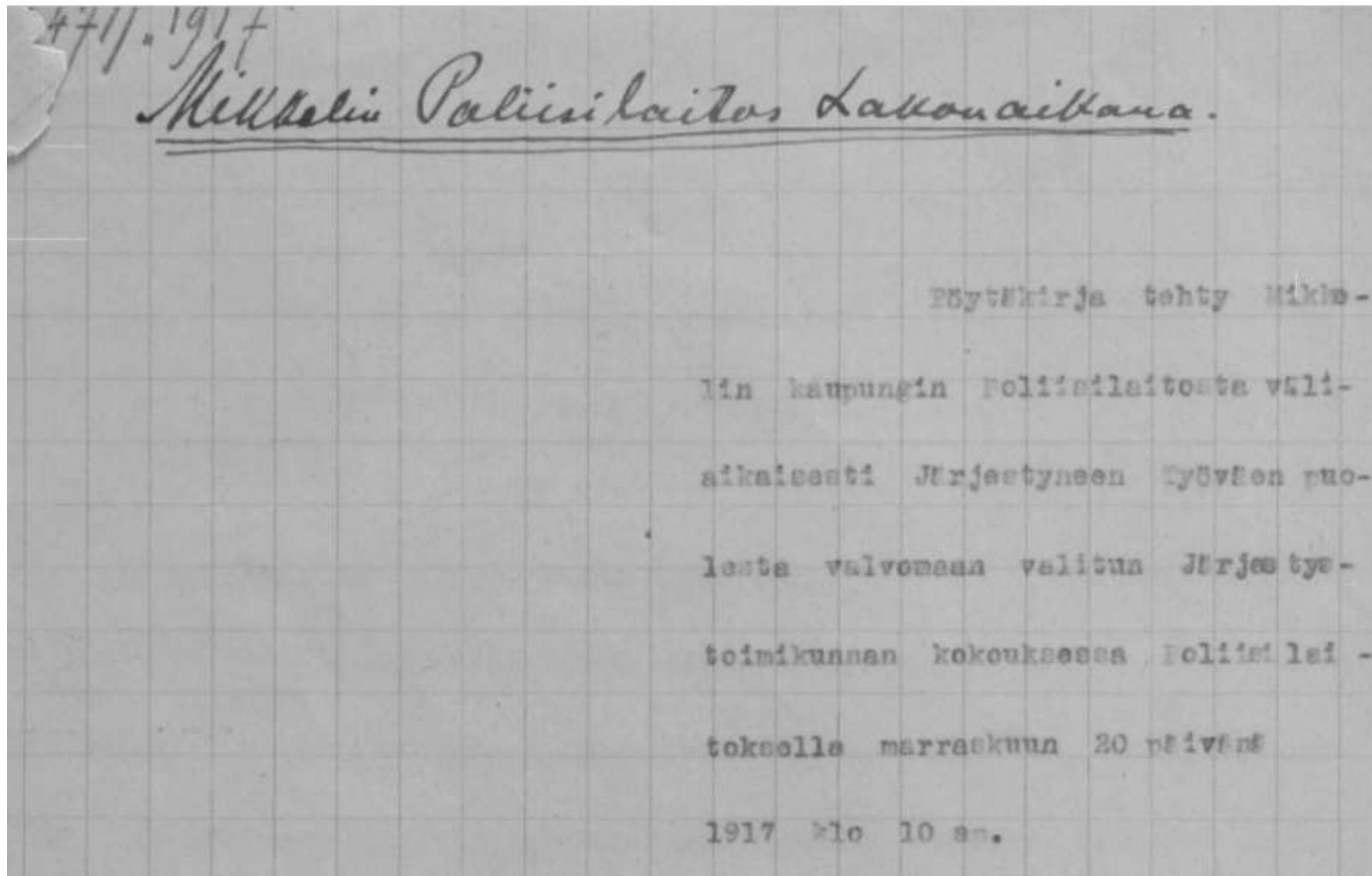| | | | | |
|---|---|---|---|---|
| Image enhancements | Image to text | Enhanced OCR & LLM Fixes | Segmentation | Keywords from text |
| Similarity analysis | Image recognition | LLM translations | Speech to Text | Summaries |
| Text to Image | Classification / tagging | Q & A | RAG | Extending information |

European Commission

# Sample image used for the presentation

# Image enhancements
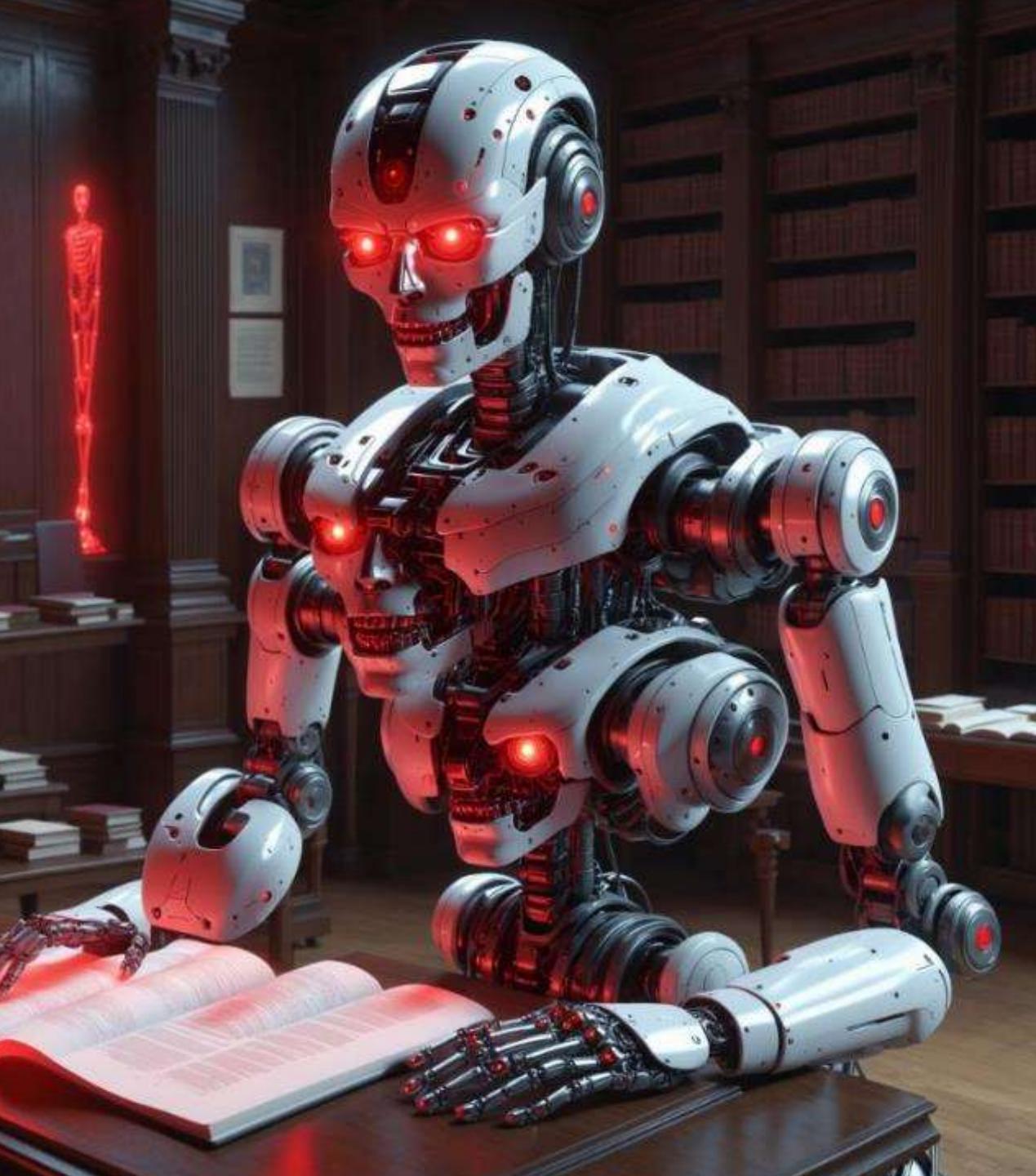
# Image to text

- You are all familiar with OCR/HTR technologies

  - Tesseract, PeroOCR, TrOCR, Abbyy fine reader, etc.

  - Based on trained AI models

  - We utilize finetuned PaddleOCR engine

- Steals parts from DLM presentation

# Enhanced OCR & LLM fixes

| Image | Original OCR | Tesseract OCR (best of enhanced) | Enhanced PaddleOCR |
|---|---|---|---|
|  | PKytrk tohty TM<br>.& ituunkin rollir ilaito.<br>ta vsii-<br>aikai eeelti Jfrjäätyneen<br>lyövf.ei^ ruo-<br>lj.-.uö iblvomeiqA yaii<br>tua J* r Jeu tyr *<br>toimikunnan<br>kflikookcacje.a . olilr4<br>lei<br>tokcolle rr ptkunn0 pr in<br>n*<br>'lo | PAY atl .<br>} Metitet: Cp leiat Le cYoq<br>Habeas alefn-tal<br><br>\| Péythyirja tehty Wikig-<br>lin \| katpungin \| potttetrat tos te<br>we1t-<br>\| aikatcestt Itnjeetyneen lydvfien<br>puo-<br>\| leste valvomaga yalitun Jirjes<br>tye-<br>toimikunnan kokouksessa<br>Polite! lei -\|<br>\| toksella marraskunn 20 ptivint<br>1917 lo 10 a. | Pöytäkirja tehty Mikle.<br>lin  kaupungin<br>Foliisilaitosta Vali<br>aikaisesti<br>Jrjestyneen<br>Työvlen puo-<br>lesta valvomaan Valitun<br>JLrjes tye-<br>toimikunnan<br>kokouksessa Foliisi lai<br>toksella<br>marraskuun  20 ptivanf<br>1917 klo10 an, |

# OCR enhacements → CER

| Model | ELKA test (4273 row pictures) | NAF1 test (3475 row pictures) | NAF2 test (3247 row pictures) | ELKA test handwritten (714 row pictures) |
|---|---|---|---|---|
| **Tesseract** | 4,6% | 2,7% | 4,4% | 79,3% |
| **Original PaddleOCR** | 6,7% | 3,9% | 6,8% | 50,2% |
| **Enhanced PaddleOCR** | 2% | 1,2% | 2,3% | 20,7% |

# LLMs in OCR corrections

- **Best OCR result**

*Pöytäkirja tehty Mikke,*
*lin  kaupungin Poliisilaitosta Vali*
*aikaisesti*
*Jrjestyneen*
*Työvlen puo-*
*lesta valvomaan Valitun*
*JLrjes tye-*
*toimikunnan*
*kokouksessa Foliisilai*
*toksella*
*marraskuun  20 ptivanl*
*1917 klo10 an,*

- **Could LLM(s) be used to fix the rest..?**
  - Copilot: Pöytäkirja tehty Mikkelin kaupungin Poliisilaitosta vali- aikaisesti järjestyneen Työväen puolesta valvomaan valitun järjestyneen työtoimikunnan kokouksessa Poliisilaitoksella marraskuun 20. päivänä 1917 klo 10 aam.
  - Gemini: Pöytäkirja tehty Mikkelin kaupungin Poliisilaitoksella aikaisesti Järjestyksen Ylläpitäjän puolesta valvomaan Valitun Järjestys-toimikunnan kokouksessa Poliisilaitoksella marraskuun 20 päivänä 1917 klo 10 a.m.
- There are great open LLMs also for Finnish & over 100 other languages

# Translations

- <mark>Huggingface:</mark> Helsinki-NLP/opus-mt-fi-en

  - Minutes made at the meeting of the elected organised working committee at the Police Department on 20 November 1917 at 10 a.m.

- <mark>CoPilot:</mark>

  - Minutes made at the Police Station of the city of Mikkeli on behalf of the temporarily organized Workers' Committee meeting at the Police Station on November 20, 1917, at 10 a.m.

- <mark>Gemini:</mark>

  - Minutes of the meeting of the organized labor committee elected to oversee the interests of the working class, held at the Mikkeli City Police Station on November 20, 1917, at 10 AM.

- <mark>Google translate:</mark>

  - The minutes were made on behalf of the labor committee organized by election from the Police Department of the city of Mikkeli to oversee the meeting of the labor committee at the Police Department on November 20, 1917 at 10 a.m.
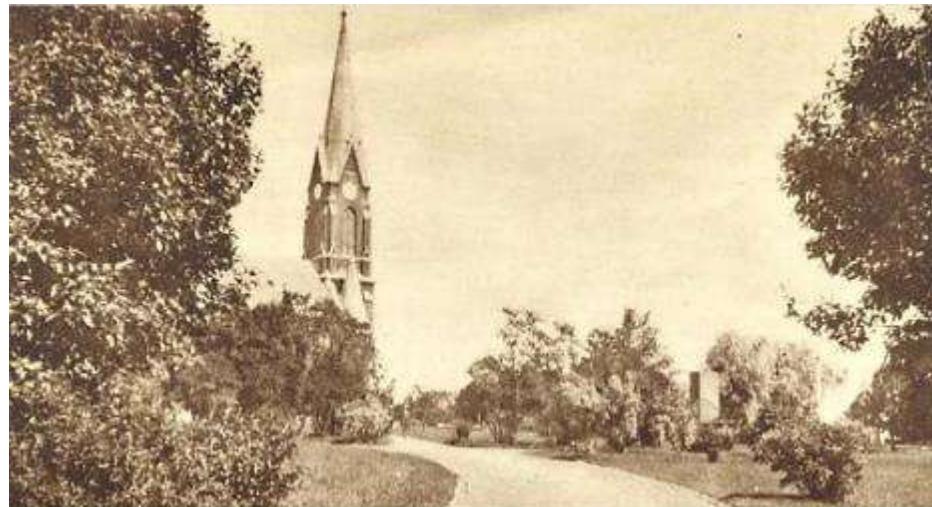
# Keywords from text

- Finnish National Library (Annif tool, https://annif.org/ ) embedded into https://arkkiivi.memorylab.fi/

| | Original OCR | PaddleOCR | LLM fixed | Eng translated |
|---|---|---|---|---|
| Index terms | tyrni, idut, suomen kieli, toimikunnat, sähkökitara, laboratoriotekniikka, terveysvaikutukset, paneelit, viljely, siirtolapuutarhat | pöytäkirjat, kokoukset, historia, järjestöt, toimikunnat, hallinto, lehtialaindeksi, poliisi (organisaatiot), kaupungit, historiikit | historia, pöytäkirjat, poliisilaitokset, poliisi (organisaatiot), historiikit, poliisihallinto, kaupungit, kokoukset, lainsäädäntö, rikokset | police (organisations), towns and cities, history, police stations, commissions (organs), police (occupations), local histories (literary works), workers, meetings (arranged events), guides (literary works) |
| Lang | fi | fi | fi | en |
| Date | | marraskuun 20 ptivanf 1917 | marraskuun 20. päivänä | November 20, 1917 |

# Summaries & classification / tagging

- In addition to keywords, LLMs support a lot more tasks

- Public demonstration: https://memorylab.fi/AIKO/cluetail-demo/

- Tested with the translated text

  - Similarity of two translations: 70.96%

  - Tone: neutral : 95.91%, positive : 2.81%, negative : 1.29%

  - Summary: Minutes made at the Police Station of the city of Mikkeli on behalf of the temporarily organized Workers' Committee.

  - Tags: Mikkeli, : 99.36%, November 20 1917, : 99.02%, Workers' Committee meeting, : 98.96% 10 a.m. : 98.69% police station, : 94.81%

# Segmentation / Image recognition & classification



The town church and surroundings

- Huggingface:

  - **microsoft/resnet-50:** church, church building 34,5%

  - **hustvl/yolos-small:** Clock 73,8%

  - **nvidia/segformer-b0-finetuned-ade-512-512:** wall, building, sky, tree, earth, plant, painting, signboard, clock

# Extending information

The town church and surroundings

beautiful place in the vicinity is **Lamposaari** island, on which there is a popular sumnner restaurant. The tourist will doubtless carry away with him a permanent memory of a visit to the famous **battlefield of Porrassalmi**. The road there leads through a lovely landscape, and arrived there the tourist sees opening before him one of the most beautiful bits of the interior, with waters glittering between high ridges and verdant fields. At the highest point of one of the ridges soars a monument, a reminder of the Bateles the Finns have had to endure in defence of their fair country. A villa on the bank of the straits has been furnished as a comfortable summer restaurant. While visiting Porrassalmi the opport unity can be taken to visit the Invalids' Home situated in the old **manor of Kyyhkylä.**

**Lamposaari** is a picturesque island located near Mikkeli, known for its natural beauty and recreational opportunities. The island features a popular summer restaurant, making it a favorite spot for both locals and tourists. The serene environment and scenic views of Lake Saimaa add to its charm.

**The Battle of Porrassalmi**, fought on June 13, 1789, during the Russo-Swedish War, is a significant historical event for Mikkeli. A Swedish force of about 750 men successfully defended against a much larger Russian force of 5,000 men. The battlefield is now a historical site with a monument commemorating the bravery and sacrifices of the Finnish soldiers. The area offers stunning views of the surrounding landscape, with glittering waters and verdant fields, making it a memorable visit for history enthusiasts and nature lovers alike.

**Kyyhkylä Manor**, established in the 1630s, has a rich history and has served various purposes over the centuries. It was originally formed from two neighboring houses and has been owned by several notable figures, including Otto Carl von Fieandt. Today, it functions as a hotel and wellness resort, offering a blend of historical charm and modern amenities. The manor's beautiful surroundings and historical significance make it a unique destination for visitors.

# Visualizing information (Stable diffusion)

Which one is real
Kyyhkylä Manor?

Is it this one

Or maybe
this one

Original

More
medieval..?

# Thank you

**Contact:**
Anssi.jaaskelainen@xamk.fi
LinkedIn
Digitalia.fi / memorylab.fi

European Commission