

Median  
museo  
ja arkisto|

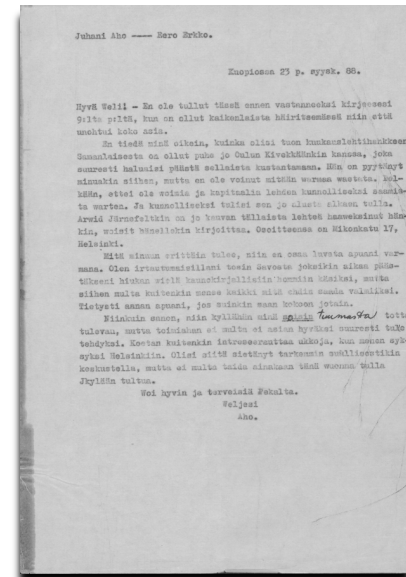
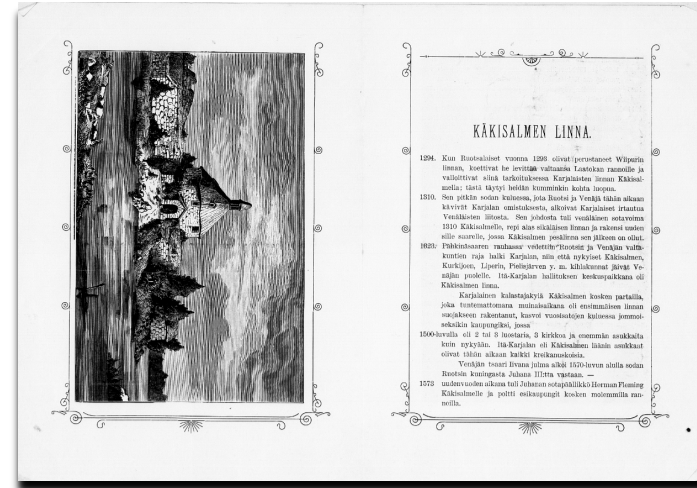
# Merkin OCR-testailut

# Tehtävänanto

- 10-15 esimerkkiä kuvista, joissa tekstiä, joka jo aiemmin tunnistettu
  - Viedään tänne: <https://memorylab.fi/AIDA/paddle-demo/>
  - Verrataan uusia OCR-tunnistuksia vanhoihin

# Testattavat aineistot

- 8 kpl digitoituja kuvia: lehtileikkeitä, koneella kirjoitettuja kirjeitä, pienpainate
- Pdf-A ja TIFF –formaatit (muutamaa JPG)
- OCR-tunnistusta tehty lähinnä lehtiaineistoihin
- Teetetty digitoitintifirmalla
- osasta irrotin vertailutekstin Adobe Acrobat Pro:lla



# Testaaminen ja tulokset

- Formaatin kanssa alkuhankaluuksia: jpg vai tiff?
  - Meillä yleisin tiedostomuoto on pdf-A
- Tulokset
  - Digitalian Paddle-versio parempi 3/8
  - Oma / aiempi versio parempi 4/8
- Erikoishuomiot
  - Värillinen sanomalehtiaineisto Paddlelle ongelmallinen
  - Paddle ei tunnistanut kaikkea tekstiä, vaikka tunnistettu oli laadukasta
  - Paddle tunnisti jopa käsinkirjoitetun tekstin
  - Adoben OCR-työkalu ei ole ihan huono 😊

# Bonustesti



SE NYT ON YKSI MERKILLINEN ASIA, ETTÄ MEIDÄN EI SALLITA PUKEA EDES KOULULAISIA MUSTIIN PAITOIHIN, VAIKKA VASEMMALLA 78 KANSANEDUSTAJAA SAA VALMISTAA KAPINAA. - NI KIRJOITTA - , ETTÄ KUN SUOMEN TALONPOJAT TULIVAT REP-PUINEEN MÄNTSÄLÄÄN SYÖMÄÄN NIIN HALLITUS TULI TYKKIEN KANSSA HEIDÄN RUOKARAUHAANSA HÄIRITSEMÄÄN. AVUN KAIKEN TÄLLÄISEN ESTÄMISEKSI ON LÄHDETTÄVÄ KANSAN RIVEISTÄ.

RIIPINEN