

SAIMAAN AMMATTIKORKEAKOULU
Tekniikka, Lappeenranta
Tietotekniikka
Viestintätekniikka

Riku Heikkilä

Käsin ladotun aikakauslehtitekstin muuttaminen digitaaliseen muotoon

TIIVISTELMÄ

Riku Heikkilä

Käsin ladotun aikakauslehtitekstin muuttaminen digitaaliseen muotoon, 49 sivua

Saimaan ammattikorkeakoulu, Lappeenranta

Tekniikka, Tietotekniikan koulutusohjelma

Viestinnän suuntautumisvaihtoehto

Opinnäytetyö 2011

Ohjaajat: Tuntiopettaja Yrjö Utti, Saimaan ammattikorkeakoulu Oy,
Jukka Seppänen, Maasotakoulu

Opinnäytetyönä tehtiin selvitys, kuinka ja millä testintunnistusohjelmalla saada toisen maailmansodan aikaiset, käsin ladotut Sotilaspaperin aikakauslehdet digitaaliseen muotoon tutkimuskäyttöön. Selvityksessä käytettiin tekstintunnistusta niin, että siihen voi tehdä sanahakuja. Opinnäytetyön asiakkaana toimi Maasotakoulun sotilaspastori Jukka Seppänen. Aineisto ja tutkimusvälineet kerättiin Internetistä.

Selvitys tehtiin saatavilla olevia ohjelmia käyttäen. Ohjelmia oli paljon ja niiden toimivuus tätä työtä ajatellen vaihteli suuresti. Ohjelmien ilmaisuus oli näennäistä, sillä kun ilmaisia ohjelmia löytyi, niin iso osa oli niin sanottuja maksullisten ohjelmien kokeiluversioita, joissa oli omat rajoituksensa liittyen useimmiten joko ohjelman käyttöikään, tallennettavien sivujen määrään tai molempiin.

Käyttökelpoisimmaksi vaihtoehdoksi tarjolla olevista ohjelmista osoittautui Nuance PDF Converter professional 7 -ohjelmalla tehty vuosikerta yhteen PDF-tiedostoon, johon voi tehdä sanahakuja. Haussa sillä on puutteita, pääasiassa harvaan kirjoitettujen nimien paikantamisessa. Ensimmäisessä tarkastelussa ohjelmasta virheitä ei löytynyt monta, mutta toisen tarkastelun jälkeen virheitä löytyi useampia. Tarjolla olevien ohjelmien huonon laadun vuoksi alkuperäissuunnitelmaan kuulunut tietokantavaihtoehto muuttui tutkimuksen edetessä yhä epätodennäköisemmäksi. Suuri syy tähän oli tekstintunnistuksen heikon laadun tuomat pulmat.

Asiasanat: PDF, tekstintunnistus, käsin ladottu, OCR

ABSTRACT

Riku Heikkilä

Transforming a hand stacked magazine in to digital form, 49 pages

Saimaa University of Applied Sciences, Lappeenranta

Technology, Degree Programme of Information Technology

Communications orientation

Bachelor's thesis, 2011

Instructors: Lecturer Yrjö Utti, Saimaa University of Applied Sciences,
Mr. Jukka Seppänen, Maasotakoulu

This Bachelor's thesis was conducted as an investigation of how to make a digitized version of a hand stacked Army priest magazine using text recognition to research use, so it would be possible to search certain words from it. The client of this thesis was army minister Jukka Seppänen. Materials and research equipment were gathered from the internet.

This thesis was executed with available software, which was plentiful, however the results varied a lot. As free as the programs were, there were a lot of so called trial versions of chargeable software which had their own limitations usually including low period of usage, maximum of saveable documents or both.

The most usable outcome was with Nuance PDF Converter professional 7, which allowed to create a volume to one PDF-file in which to make word searches. It still had flaws, mainly in searching infrequently written words. Because of the flimsy quality of the results from the character recognition, mainly in transforming the file in to text, the database option which was included in the original plan, had eventually to be discarded.

Keywords: PDF, volume, hand stacked, OCR

SISÄLTÖ

1	Johdanto	8
2	Tekstintunnistusohjelman valintaprosessi	9
3	Testauskohteen esittely.....	10
4	Digitointityön suoritus	13
4.1	Asiakirjojen kuvaustapa.....	13
4.2	Skannaus.....	14
5	Testatut ohjelmat.....	15
5.1	Selaimessa toimivat OCR-ohjelmat	15
5.1.1	Free OCR -ohjelma.....	19
5.1.2	Free Online OCR -ohjelma.....	21
5.1.3	i2OCR -ohjelma	22
5.1.4	New OCR -ohjelma	24
5.1.5	OCR Convert -ohjelma.....	26
5.1.6	OCR Now! -ohjelma	27
5.1.7	OCR online -ohjelma.....	28
5.1.8	OCR terminal -ohjelma.....	29
5.2	Koneelle asennettavia OCR-ohjelmia.....	31
5.2.5	Abby Finereader 10 -ohjelma	34
5.2.6	Adobe Acrobat pro 10 -ohjelma	35
5.2.7	Autobahn DX -ohjelma	36
5.2.8	Free OCR 3.0 -ohjelma	38
5.2.9	Smart OCR -ohjelma.....	39
5.2.10	PDF Converter Professional 7 -ohjelma.....	40
6	Tulosten esittely	41
7	Yhteenveto ja päätelmät.....	45
	Kuvat.....	47
	Taulukot.....	48
	Lähteet.....	49

TERMIT JA LYHENTEET

BMP	Bittikarttakuva (Bit Map Picture), joka muodostuu pikseleistä eli kuvapisteistä.
Digikamera	Tunnetaan myös nimellä digitaalikamera. Se on kamera, joka tallentaa kuvattavan kohteen digitaalisesti.
Digitalisointi	Dokumentin siirtäminen digitaaliseen muotoon.
DOC	(Document) Tekstiedostomuoto, joka on käytössä Microsoftin Word-ohjelmassa.
DOS	Disk Operating System eli levykäyttöjärjestelmä on yhden käyttäjän komentorivipohjainen ei-moniajava käyttöjärjestelmä.
GIF	Graphic Interchange Format on bittikarttagrafiikan tallennusformaatti, joka käyttää häviötöntä pakkausta.
HTML	Hypertext Markup Language eli hypertekstin merkintäkieli on avoimen standardin kuvauskieli, jolla voi kuvata tekstiä, jossa on hyperlinkkejä, eli hypertekstiä.
JPEG	Joint Photographic Experts Group on bittikarttagrafiikan tallennusformaatti, joka käyttää häviöllistä pakkausta.
Krediitti	Credit on nettisivustoilla käytettävää valuuttaa.
MB	Megabyte eli megatavu. Tavu on tiedon tallennuskapasiteetin mittayksikkö tietokoneessa.
MHT	Tiedostotyypistä käytetään myös lyhennettä MHTML (MIME HTML), joka on web-sivun arkistformaatti. Tätä käytetään yhdistämään kuvia, animaatioita ja äänitiedostoja HTML-koodin kanssa yhdeksi tiedostoksi.
OCR	Optical Character Recognition on teknologia, jonka avulla tunnistetaan koneellisesti tai käsin tuotettua tekstiä sähköiseen muokattavaan muotoon.
PBM/PGM/PPM	Portable bitmap format, Portable greymap format ja Portable pixmap format ovat avoimen lähdekoodin netpbm-projektin tuottamia kuvaformaatteja, jotka luotiin helposti siirrettäväksi eri ohjelmistojen välillä.

PCX	Personal Computer eXchange on kuvaformaatti, joka oli ensimmäinen hyväksytty DOS-kuvastandardi (Disk Operating system), mutta on vähentynyt käytöstä.
PDF	Portable Document Format on Adoben kehittämä PostScript-kieleen pohjautuva ohjelmistoriippumaton, siirrettävä tiedostomuoto, jota käytetään pääasiallisesti sähköiseen julkaisemiseen, tulostamiseen ja painamiseen.
PNG	Portable Network Graphics on bittikarttagrafiikan tallennusformaatti, joka käyttää häviötöntä pakkausta. Luotiin GIF:n (Graphic Interchange Format) korvaajaksi.
PostScript	Tietokoneiden sivunkuvauskieli, jolla kuvataan tulostimille ja muille laitteille tulostettavien dokumenttien ulkoasu.
PPT	Power Point Table on Microsoftin PowerPoint-ohjelman käyttämä tiedostomuoto.
RTF	Rich Text Format on muotoillun tekstin tallennusmuoto joka, toisin kuin tavallinen teksti (.txt), tallentaa tiedostoihin myös fontit, fonttikoot, kursivoinnit, lihavoinnit sekä alleviivaukset.
Selain	Ohjelma, jolla käyttäjä voi katsella ja lähettää kuvia, tekstiä ja muuta, mitä www-sivuilta löytyy.
Skannauskynä	Kynä, jolla voi skannata sanoja rivi kerrallaan.
Skanneri	Laite, joka muuttaa kaksiulotteisen kuvan kuvatiedostoksi.
Tekstin-tunnistus-ohjelma	Ohjelma joka tunnistaa kuvasta merkit ja tulostaa ne määrättyyn tiedostomuotoon.
TIFF	Tagged image file format tai Tag Image File Format, joista voidaan käyttää myös lyhennystä TIF. TIFF on kuvien tallennukseen käytetty tiedostomuoto, joka tukee niin kuvankäsittely- kuin tekstinkäsittelyohjelmia.
TWAIN	Technology Without An Interesting Name on kuvankäsittelyssä standardirajapinta, joka tehtiin varmistamaan, että eri laitteet ja ohjelmistot ovat yhteensopivia. Standardi mahdollistaa skannaamisen suoraan kuvankäsittelyohjelmasta.
TXT	Text file eli tekstitiedosto on tiedostomuoto, joka sisältää vain tekstiä.

Validointi	Prosessi, jolla varmistetaan, täyttääkö jokin asia tietyt vaatimukset.
WORD	Microsoftin kehittämä tekstinkäsittelyohjelma.
WWW	World Wide Web on Internet-verkossa toimiva hypertextijärjestelmä.
XLS	Excel Spreadsheet on Microsoftin Excel- taulukkolaskentaohjelman käyttämä tiedostomuoto.
XML	eXtensible Markup Language on merkintäkieli, jolla tiedon merkityksen voi kuvata tiedon sekaan. Käytetään tiedonvälitykseen järjestelmien välillä ja formaattina dokumenttien tallentamiseen.

1 Johdanto

Yleensä digitoinnista eli digitalisoinnista puhuttaessa viitataan valokuvien, dioiden tai mikrofilmien digitalisointiin. Dokumenttien muuttaminen digitaaliseen muotoon on aivan yhtä lailla digitalisointia, vaikka useammin tekstidokumenttien digitalisoinnista käytetään nimitystä tekstintunnistus tai OCR.

Opinnäytetyön tavoitteena on tutkia keinoja muuttaa toisen maailmansodan aikainen Sotilaspaperin aikakauslehti digitaaliseen muotoon asiasanojen hakua varten. Tarkoituksena on löytää ilmainen OCR- (Optical character recognition) eli tekstintunnistusohjelma, joka tunnistaa merkit mahdollisimman pienellä virhemarginaalilla. Ohjelmasta tehdään joko alkuperäisen tekstin siirto määritettyyn tietokantaan tai yksinkertaisemmin muutetaan tiedosto PDF:ksi (Portable Document Format), johon voi tehdä sanahakuja. PDF on sähköiseen julkaisemiseen, tulostamiseen ja painamiseen käytetty ohjelmistoriippumaton ja siirrettävä tiedostomuoto. Työn taustalla on aikakauslehdien digitaalisen version tarve, koska digitaalisessa muodossa tieto säilyy ja tiedon etsiminen sisällöstä helpottuu.

Ohjelmien etsiminen tapahtuu Googlen avulla. Ohjelmat testataan tekstintunnistuksen jälkeen tutkimalla virheiden määrää, jonka jälkeen virheiden määriä vertaillaan ohjelmien kesken. Sen jälkeen valitaan ohjelma, jolla vuosikerrasta pystytään tekemään digitaalinen versio ja josta löytyy vähiten virheitä.

Tässä työssä käsitellään vain tekstintunnistusohjelmien ilmaisversioita, sillä maksulliset versiot voivat hinnaltaan kivuta useampaan sataan euroon eikä sellaiseen ollut tällä kertaa mahdollisuutta. Ensin tutkitaan selainpohjaisia OCR-ohjelmia, jotka eivät tarvitse lainkaan työasemaan asentamista toimiakseen, ja tämän jälkeen perehdytään asennettaviin ohjelmiin.

2 Tekstintunnistusohjelman valintaprosessi

Testausmateriaalin määrittäminen aloitetaan määrittämällä testausmateriaali eli valitaan kuvatiedosto, joka on sisällöltään monipuolista. Sisällöstä on löydettävä kaikki lähdemateriaalista löytyvät tekstityypit esimerkiksi eri tavoin painetut kirjaimet, harvaan kirjoitetut sanat ja eri tasoille painetut sanat. Vaihtelevuuden ja kattavuuden haluamiseksi testimateriaaleja voi olla pari lisää.

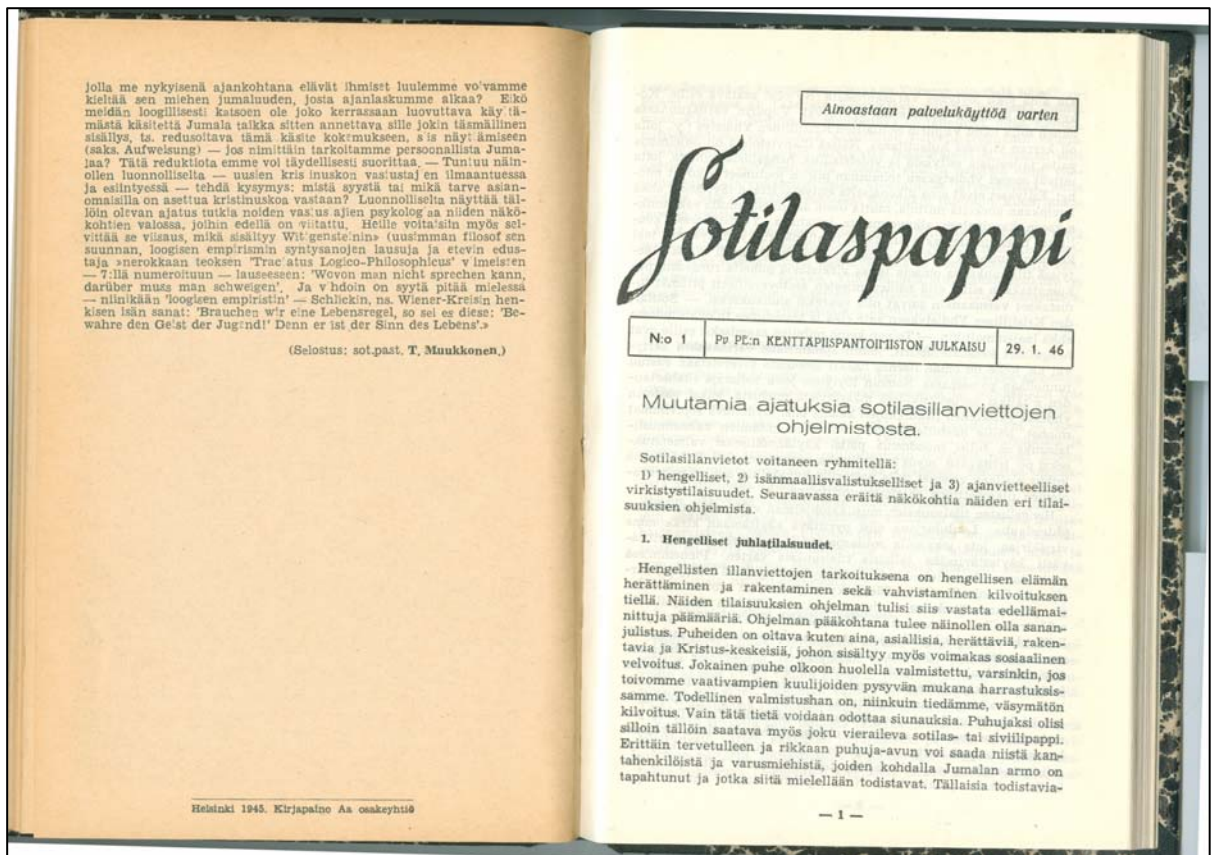
Internetistä etsitään kaikki tekstintunnistukseen viittaavat ohjelmat. Tehdään alustava tutkimus, onko ohjelmia mahdollista testata.

Ohjelmat testataan, kun testattavat ohjelmat on löydetty. Testattujen ohjelmien tuloksia verrataan keskenään. Pienimmän virhemäärän tulostanut ohjelma valitaan. Seuraavaksi materiaali käsitellään valitulla ohjelmalla.

Skannausprosessi ohjelman ollessa tuotannossa: ensin kohde määritetään, jonka jälkeen lähdemateriaali skannataan. Materiaali käsitellään tekstintunnistusohjelmassa, jonka jälkeen tulokset tarkistetaan ja verrataan lähdemateriaaliin. Lopuksi saadut tulokset siirretään Internetiin.

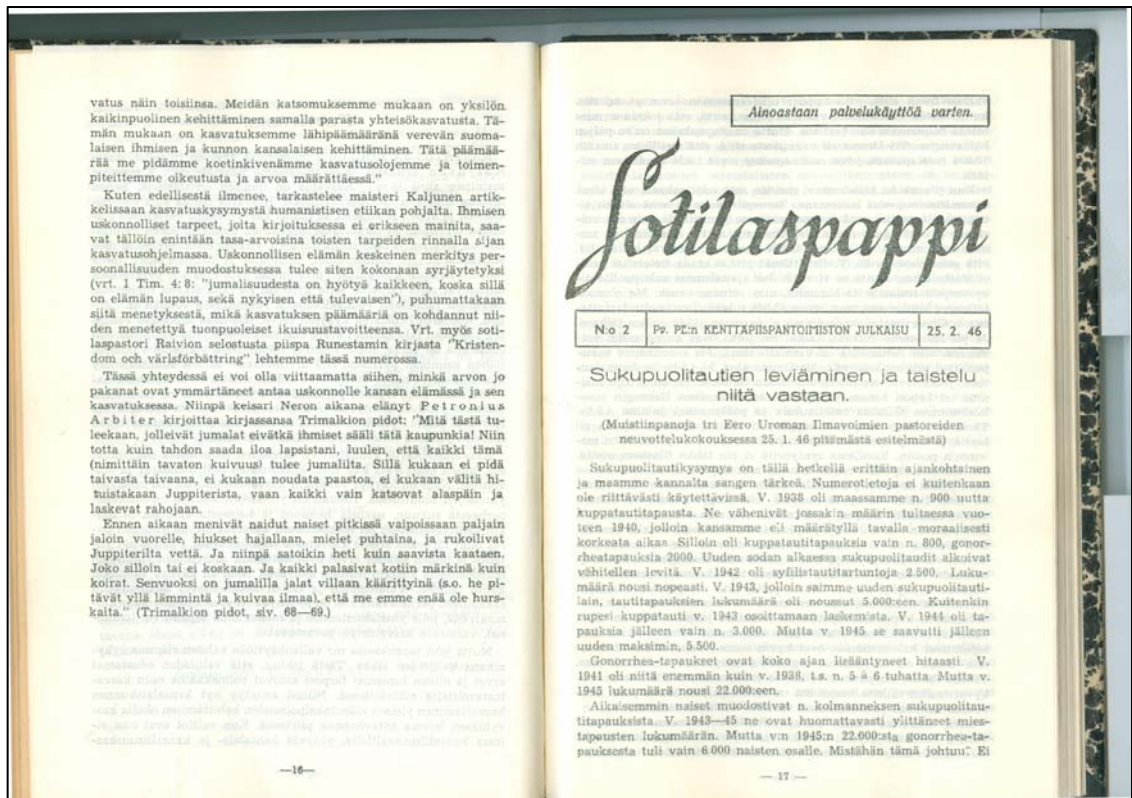
3 Testauskohteen esittely

Projektin kohteena on Sotilaspappi-niminen kenttäpiispan toimiston käsin ladottu aikakauslehti (kuva 1). Käsin ladottu tarkoittaa sitä, että kirjaimet ovat painoa varten käsin asetetut, jonka takia teksti ei saata olla joka rivillä samalla tasolla. Testattavana oli 46. vuosikerta, joka sisältää 114 sivua, jotka skannattiin eli muutettiin kaksiulotteinen kuva kuvatiedostoksi. Skannauksen tuloksena saatiin 59 kuvatiedostoa eli kuva aukeamaa kohti. Lehti sisältää isoja ja paksuja sanoja, harvaan kirjoitettuja sanoja ja erikoismerkkejä kuten heittomerkkejä, sulkumerkkejä ja kaksoispisteitä.

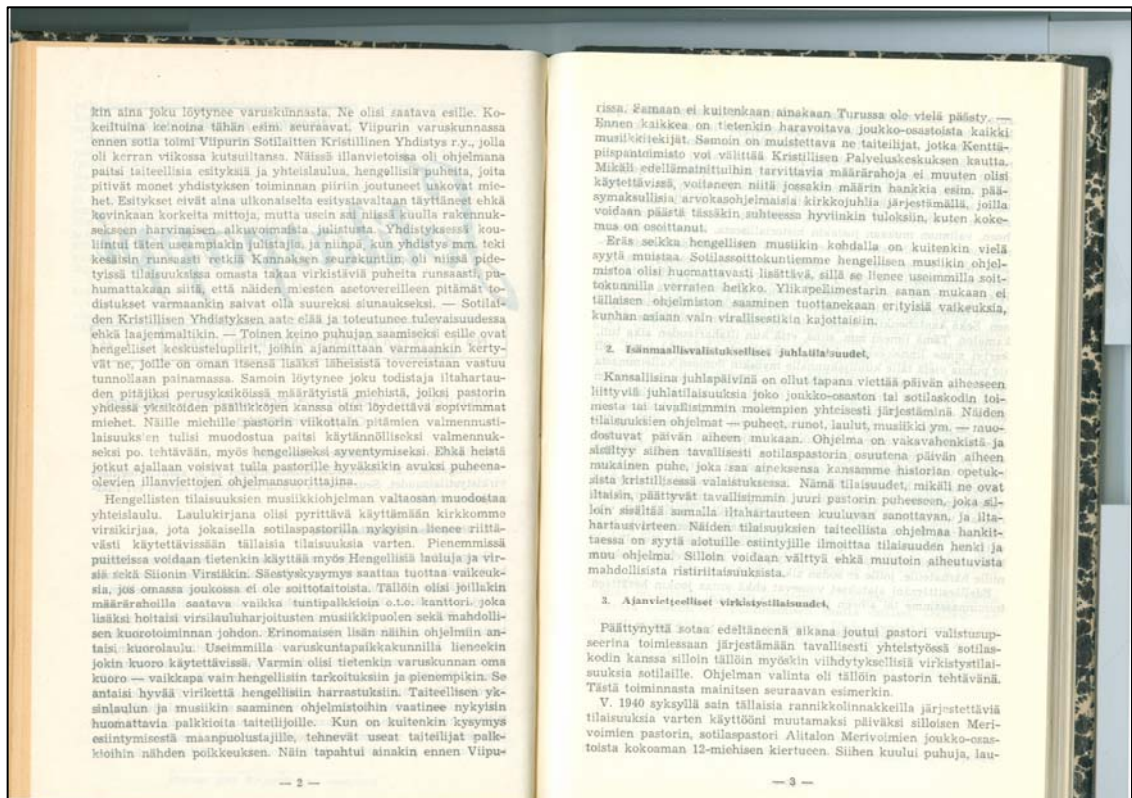


Kuva 1. Sotilaspappi-aikakausilehden ensimmäinen sivu

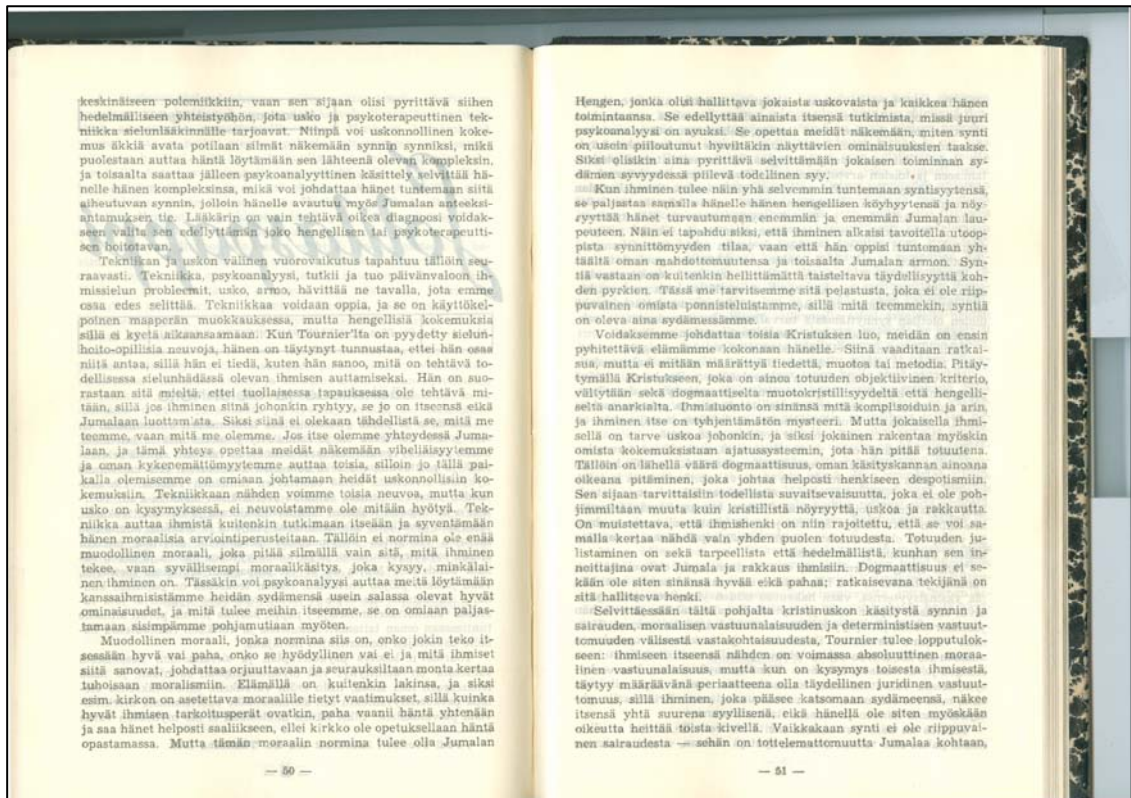
Eri sivujen välillä oli ladonnan tasaisuuden vaihteluja, ja se asetti tekstintunnistukselle haasteita. Testauskohteina käytettiin viittä (5) eri aukeamaa, jotta saatiin kattava ja edustava valikoima tekstistä. Testausaineistona 1 oli testatun vuosikerran ensimmäinen aukeama (Kuva 1). Eniten käytössä oli testausaineisto 2 (Kuva 2). Käytössä olivat myös Testausaineistot 3 (Kuva 3), 4 (Kuva 4) ja 5 (kuva 5). Kaikki testiaineistot ovat vuosikerrasta 46.



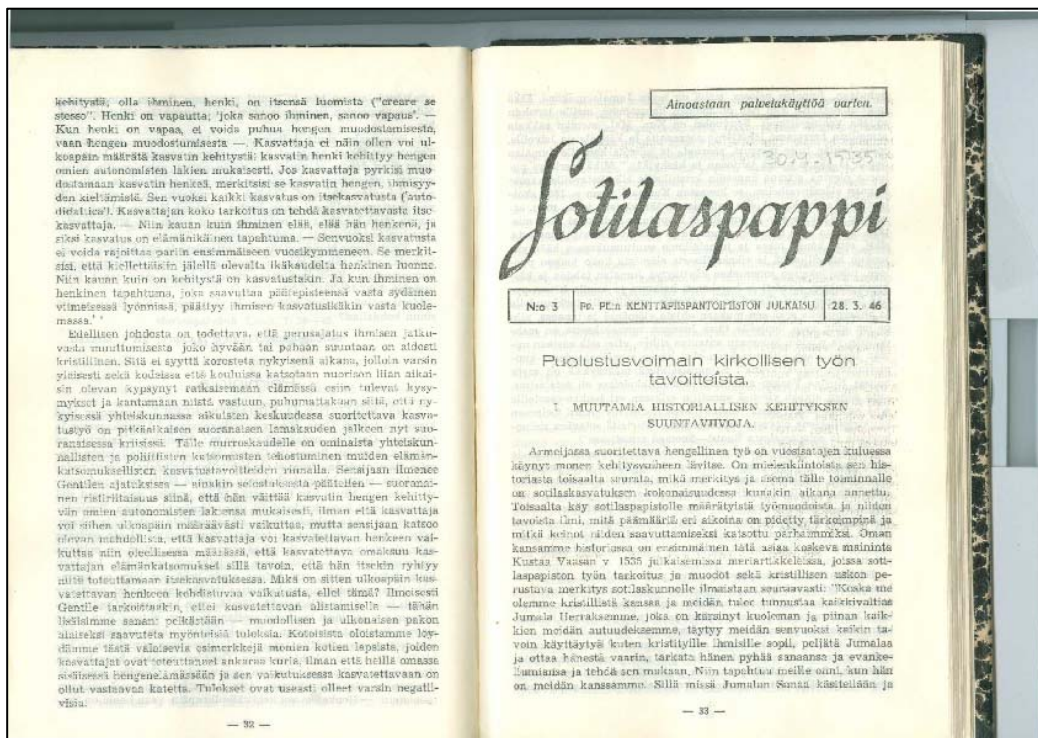
Kuva 2. Testausaineisto 2:n esittely (sivut 16 ja 17, jolla on luvun kaksi alku)



Kuva 3. Testauskohteen 3 esittely (sivut 2 ja 3 luvusta 1)



Kuva 4. Testiaineisto 4:n esittely (sivut 50 ja 51 luvusta 4-5)



Kuva 5. Testiaineisto 5:n esittely (sivut 32 ja 33, jolla on luvun kolme alku)

Tasaisten lopputulosten saavuttamiseksi testiaineisto 2 oli eniten testattu kohde, mutta vaihteluvuoden ja kattavuuden vuoksi testiaineisto 1, 3, 4 ja 5 olivat

myös käytössä. Kuten kuvista on nähtävissä, oman haasteen tekee myös sivujen kääntöpuolen tekstin näkyminen paperin läpi.

4 Digitointityön suoritus

Digitointi- eli digitalisointityö suoritettiin skannaamalla nidottu lehti, jonka jälkeen skannatut kuvat käsiteltiin testattavilla tekstintunnistusohjelmilla. Ohjelmia käytettiin kahdenlaisia. Ensin testattiin selainpohjaisia ohjelmia, ja niiden jälkeen työasemalle asennettavia ohjelmia.

4.1 Asiakirjojen kuvaustapa

Asiakirjojen kuvaamiseen löytyy kaksi vaihtoehtoista tapaa: skanneri ja digikamera, eli kamera, joka tallentaa kuvattavan kohteen digitaaliseksi. Tavallisesti skannerit eivät välttämättä sovellu paksujen ja nidottujen teosten skannaamiseen, koska se voi vahingoittaa aineistoa. Sivut voivat irrota tai teksti suttaantua. Digikameralla kuvattuna puolestaan mittasuhteet ja linjat voivat vääristyä. Skannauskynää ei kokeiltu, koska aineistoa oli satoja sivuja ja skannauskynällä skannataan vain rivi kerrallaan.

"Kamera on skanneria sopivampi laite digitointiin silloin, kun esimerkiksi kirjaa ei voida avata levälleen tai digitoitava materiaali on haurasta ja vaarassa vahingoittua. Kameraa käytettäessä ei laitteen ja kopioitavan materiaalin välillä ole fyysistä kontaktia, joka saattaisi vahingoittaa alkuperäistä aineistoa."

"Valokuvaaminen digitaalikameralla saattaa olla ainoa mahdollisuus digitoida materiaalia, jota ei voida muilla keinoin käsitellä esimerkiksi aineiston vahingoittumattomuuden varmistamiseksi. Sen käytössä on omat riskinsä. Esimerkiksi kamerasobjektiivi saattaa vääristää kohteen mittasuhteita ja linjoja, tai tiedoston resoluutio ja tiedostomuoto saattavat poiketa siitä, mikä on organisaation standardeissa."

"Digitaalikameran käyttökelpoisuus digitalisoinnissa onkin harkittava huolella."

/1/

4.2 Skannaus

Skannaukseen valittiin HP Scanjet G4010 -valokuvaskanneri. Asetuksissa täytyi muistaa asettaa TWAIN (Technology Without An Interesting Name) päälle ohjeiden mukaisesti. TWAIN on standardi, joka varmistaa eri laitteiden ja ohjelmistojen yhteensopivuuden sekä mahdollistaa skannaamisen kuvankäsittelyohjelmasta. Skannaus tehtiin väriskannauksella.

Yksi aukeama skannattiin kerrallaan. Skannatut tiedostot tallennettiin tiedostomuotoon TIFF (Tagged image file format), joka tunnetaan myös lyhenteellä TIF. TIFF on kuvien tallennukseen käytetty tiedostomuoto. Tiedostot nimettiin sivunumeroidensa mukaan kolmella numerolla, esimerkiksi 001, 014 ja 127, etteivät sivut tulostu kuvia yhdistettäessä väärässä järjestyksessä. Tiedoston nimeen merkittiin myös vuosikerta sekä luku, jotta tiedostot eivät olisi menneet sekaisin. Kun tiedosto oli käsitelty OCR-ohjelmalla, käsitellyn kuvan nimeen lisättiin ohjelman nimi, jotta testatut tiedostot eivät menisi sekaisin eri ohjelmien välillä. 144-sivuisesta vuosikerrasta 46 kuvatiedostoa tulostui yhteensä 58 kappaletta

5 Testatut ohjelmat

Työssä testattiin kahdenlaisia ohjelmia: niitä, jotka toimivat Internet-selaimesta käsin, ja niitä, jotka vaativat asennuksen työasemalle. Tarkoituksena oli muuttaa skannattu kuvatiedosto joko DOC-tiedostoksi (Document) - joka on Microsoftin käyttämä tekstitiedostomuoto - tai PDF-tiedostoksi tai molemmiksi ja katsoa, miltä lopputulos näyttää. Muutoksen valinta riippui ohjelmasta ja siitä, mitä mahdollisuuksia ohjelma antaa tunnistusformaatin suhteen.

5.1 Selaimessa toimivat OCR-ohjelmat

Ensimmäisenä testattiin, kuinka Internet-selaimesta käsin toimivat tekstintunnistusohjelmat toimivat. Selain on ohjelma, jolla käyttäjä voi katsella ja lähettää kuvia, tekstiä ja muuta, mitä WWW-sivuilta (World Wide Web eli maailmanlaajuinen verkko) löytyy. Testattavia ohjelmia kokeiltiin kahdeksan (8) kappaletta ja tutkittiin, kuinka hyvä lopputulos niillä saadaan. Ohjelmia oli kahdenlaisia: niitä, jotka toimivat rajoituksetta, ja niitä, jotka muuttuivat 2 – 20 käsitellyn tiedoston jälkeen maksullisiksi. Ohjelmissa oli eroavaisuuksia tuettujen kielten suhteen (Taulukko 1). Jotkut ohjelmat tukivat monia kieliä, ja jotkut tunnistivat vain englanninkielestä löytyviä merkkejä ja sanoja.

Ohjelmat erosivat siinä, mitä kuvaformaattia ne tukevat (Taulukko 3) ja mihin tiedostomuotoon käsitelty tiedosto tulostuu (Taulukko 4). Tuettuja kuvaformaatteja oli runsaasti: PDF, GIF (Graphic Interchange Format), TIFF, JPEG (Photographic Experts Group) ja PNG (Portable Network Graphics), jotka ovat bittikarttagrafiikan tallennusformaatteja. BMP (Bit Map Picture) on bittikarttakuva. PGM (Portable Greymap Format) ja PPM (Portable Pixmap Format) ovat avoimen lähdekoodin kuvaformaatteja. PCX (Personal Computer eXchange) on kuvaformaatti, joka oli ensimmäinen hyväksytty DOS-kuvastandardi (Disc Operating System). DOS on yhden käyttäjän komentorivipohjainen ei-moniajava käyttöjärjestelmä.

Tiedostomuotoja, jotka voidaan valita käsittelyn jälkeiseksi formaatiksi, oli muutamia: DOC, PDF, TXT (Text file) eli tekstitiedosto, RTF (Rich Text Format), joka on muotoillun tekstin tallennusmuoto sekä XML (eXtensible Markup Lan-

guage), joka on dokumenttien tallennusformaatti ja jota käytetään järjestelmien väliseen tiedonvälitykseen. Lopuksi on XLS (Excel Spreadsheet), joka on Microsoftin Excel-taulukkolaskentaohjelman käyttämä tiedostomuoto.

Ohjelma	suomi	englanti	ruotsi	muut kielet yht.
Free OCR	x	x	x	26
Free Online OCR		x		
i2OCR	x	x	x	30
New ocr	x	x	x	26
OCR Convert		x		4
OCR Now!	x	x	x	13
Ocr online	x	x	x	150
OCR terminal	x	x	x	150

Taulukko 1. Selainpohjaisten tekstintunnistusohjelmien tunnistamat kielet

Käsittelyissä käytettiin neljää (4) eri testiaineistoa (Taulukko 2), jotta saataisiin erilaisia näkökulmia ja vaihtelevuutta. Kielten tukeminen oli ohjelmissa vaihtelevaa ja OCR Online – ohjelmassa joutui asetuksista lisäämään suomen kieliin, koska listassa on aluksi vain 15 kieltä, joita pystyy asetuksista valitsemalla helposti muokkaamaan. Yhteensä kielivaihtoehtoja on 153.

Ohjelma	Testiaineisto 1	Testiaineisto 2	Testiaineisto 3	Testiaineisto 4
Free OCR	x			
Free Online OCR		x		
i2OCR		x		
New ocr	x	x		
OCR Convert		x		
OCR Now!		x		
Ocr online			x	
OCR terminal				x

Taulukko 2. Ohjelmissa käytetyt testiaineistot

Ohjelmat	PDF	GIF	TIFF	JPEG	BMP	PNG	PGM	PPM	PCX
Free OCR	x	x	x	x	x				
Free Online OCR	x	x	x	x	x	x			
i2OCR		x	x	x	x	x	x	x	
New ocr		x	x	x	x	x			
OCR Convert	x	x		x	x				
OCR Now!	x	x	x	x		x			
Ocr online	x	x	x	x		x			
OCR terminal		x	x	x	x	x			x

Taulukko 3. Selainpohjaisten tekstintunnistusohjelmien käsiteltävät kuvaformaattit

Ohjelma	DOC	PDF	TXT	RTF	XML	XLS
Free OCR	x	x	x	x		
Free Online OCR	x	x	x	x		
i2OCR	x	x		x	x	
New ocr			x			
OCR Convert			x			
OCR Now!	x				x	x
Ocr online	x	x	x	x		
OCR terminal	x	x	x	x		x

Taulukko 4. Selainpohjaisten tekstintunnistusohjelmien tulostettavat tiedostomuodot

Ohjelmien rajoitukset (Taulukko 5) vaihtelivat ilmaisuuden ja hintatason välillä.

Ohjelma	Ei rajoituksia	Hinta	Kokeilu-aika	Tiedostorajoitus	Ladattavan tiedoston koko (MB)
Free OCR					2
Free Online OCR	x				
i2OCR	x				
New OCR					5
OCR Convert	x				
OCR Now!		1,99 puntaa		2	
Ocr online		7,99 dollaria		10/viikko	
OCR terminal		5 dollaria		20/kuukausi	

Taulukko 5. Selainpohjaisten tekstintunnistusohjelmien rajoitukset

Hinta on minimihinta, joka täytyy maksaa, kun tiedostorajoitus täyttyy. Minimihinta tarkoittaa sitä hintaa, jolla saa ostettua krediittejä (Credit), eli sivustoilla käytettävää valuuttaa. Ostettavien krediittien tai sivujen määrä vaihteli. OCR Now! -ohjelmassa 1,99 punnalla sai 10 krediittiä, OCR Online -ohjelmassa 7,99 dollarilla sai 200 sivua, joka tarkoittaa 200:aa kuvatiedostoa. OCR Terminal -ohjelmassa 50 sivua maksoi 5 dollaria.

5.1.1 Free OCR -ohjelma

Tuetut kuvaformaattit näkyvät Taulukossa 3. Kieleksi voi valita suomen, englannin ja ruotsin lisäksi 26 kieltä, joten myös ä- ja ö-kirjaimet löytyvät. Käytön validoimiseksi, eli sen varmistamiseksi, täyttääkö jokin asia tietyt vaatimukset, täytyy vain kirjoittaa ruudulla näkyvä merkkijono.

Toisin kuin moni muu ohjelma, Free OCR /2/ ei siirrä tekstiä kuvineen Wordiin (Microsoftin kehittämä tekstinkäsittelyohjelma), vaan jättää tekstin ilman kuvia selaimen vapaaseen käyttöön (Kuva 6), josta ne voi siirtää haluamaansa ohjelmaan.



Kuva 6. Free OCR -ohjelman tekstintunnistusnäyte

Seuraavanlaisesti tulostui käännös (Kuva 7) testiaineistosta 1 (Taulukko 2); Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai

tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat sanat, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisella merkityjä sanoja löytyi 20 kappaletta. Keltaisella merkityjä sanoja löytyi 6 kappaletta.

N Ainoastaan palveluräytiö varten I

Q z I

^ I 3

Pv Plizn KENUÄPIISPANTOIHISTON JULKAISU -9. I. 46

Muutamia ajatuksia sotilasillanviettojen ohjelmistosta.

Sotilasillanvietot voitaneen ryhmitellä:

1) hengelliset, 2) isä-miaallisval.istuksell:iset ja 8) aj-anvietteelliset virkistystilaisuudet. Seuraavassa eräitä niälkäkohtia näiden eri tilaisuuksien ohjelmista.

1. Hengelliset ihlivilaisuudet.

Hengellisten illanviettojen tarkoituksena on hengellisen elämän herättäminen ja rakentaminen sekä vahvistaminen kilvoituksen tiellä. Näiden tilaisuuksien ohjelman tulisi siis vastata edellämäinittuja päämääriä. Ohjelman pääkohtana tulee näinollen olla sanan-julistus. Puheiden on oltava -kuten aina, asiallisia, heräfläviä, rakentavia ja Kristus-keskeisiä? johon sisältyy myös voimakas sosiaalinen velvoitus. Jokainen gpuhe olkoon huolella valmistettu, varsinkin, jos toivomme vaativimpien kuulijoiden pysyvän mukana iharrastuksis-samme. Todellinen valmistushan on, tiedämme, väsymätön kilvoitus. Vain tietä voidaan odottaa siimaitksia. Puhujaksi olisi silloin tällöin saatava myös joku vieraileva sotilas- tai siviilipappi. Erittäin tervetulleen ja rikkaan puhuja-avun voi saada niistä kani-tahenkilöistä ja varusmiehistä-, joiden kohdalla Jumalan armo on tapahtunut ja jotka mielellään todistavat. Tällaisia todistavia-

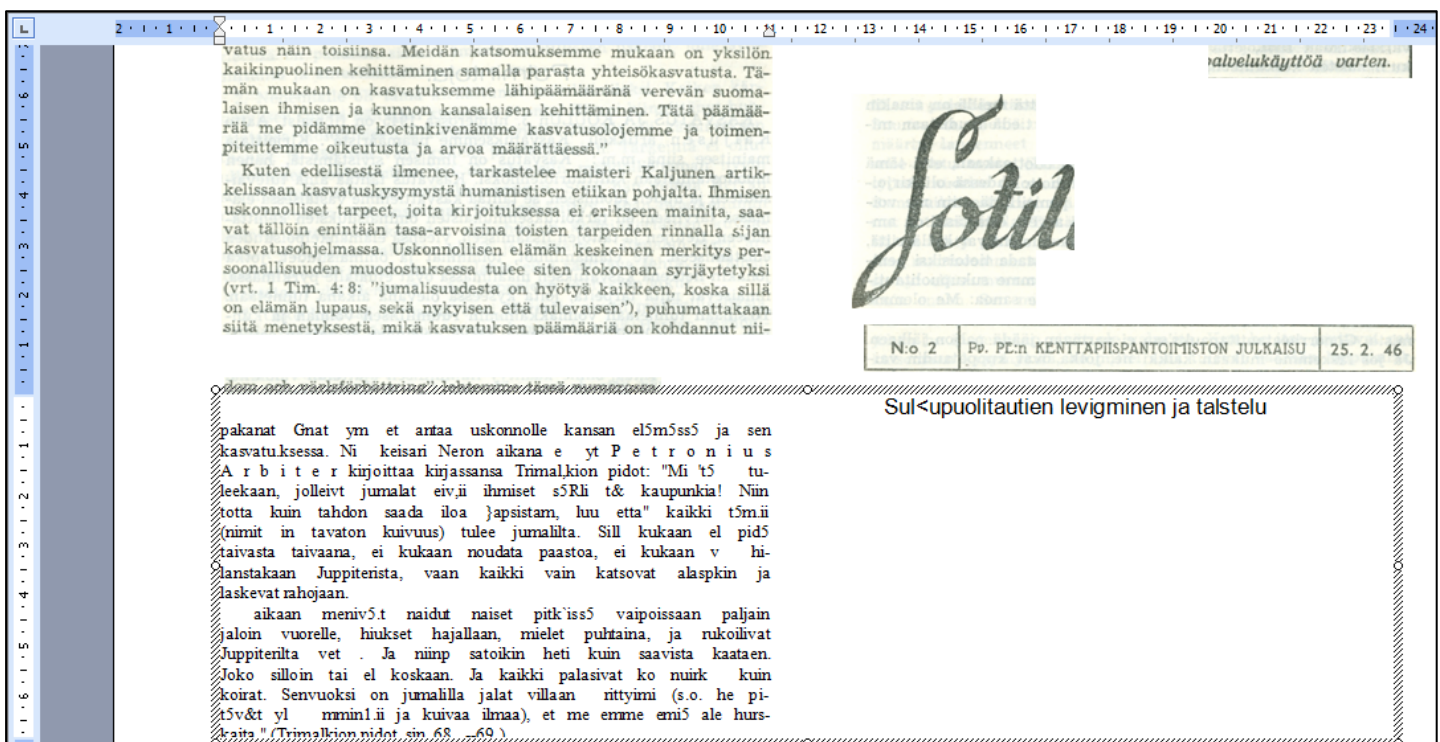
_ 1 _ V

Kuva 7. Free OCR -ohjelman tekstintunnistusnäyte testiaineistosta 1.

Lopputuloksessa oli virheitä eli tunnistamattomia merkkejä niin paljon, että ohjelman käyttäminen projektissa olisi ollut turhaa, joten virheitä ei korjattu suuren työmäärän takia.

5.1.2 Free Online OCR -ohjelma

Free online OCR /3/ on yksinkertainen suoraan selaimesta toimiva tekstintunnistusohjelma. TIFF-kuvan muuntamisen jälkeen huomattiin, että suurin osa pienistä L-kirjaimista on muuttunut f-kirjaimiksi. PDF:n muutoksessa osa tekstistä tunnistettiin, mutta virheitä löytyi, kun taas suurin osa tulostui kuvana. Käsiteltäessä Jpeg-kuva Word-tiedostoksi ohjelma tunnisti osan (kuva 8) tekstistä, mutta lähes joka toinen sana oli epäonnistunut, joten sillä ei ole mitään merkitystä, ettei ohjelma tunnistanut ä- ja ö-kirjaimia.



Kuva 8. Free online ocr -ohjelman tekstintunnistusnäyte

Free Online OCR -ohjelman testissä käytettiin testiaineistoa 2 (Taulukko 2, s. 17). Tulos oli tekstintunnistuksen kannalta heikko, sillä virheitä oli paljon. Niiden korjaaminen testissä olisi vaatinut lähes koko tekstin kirjoittamisen uudelleen.

5.1.3 i2OCR -ohjelma

Englannin ja suomen lisäksi i2OCR-ohjelmasta löytyy 31 vaihtoehtoista kieltä. Ohjelma /4/ tukee useita kuvaformaatteja (Taulukko 3, s.17). Lopputuloksessa (Kuva 9) oli yllättävän vähän virheitä, mutta silti liikaa.

Seuraavanlaisesti tulostui käännös testimateriaalista 2 (Taulukko 2, s.17). Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat sanat, joissa ohjelma on tunnistanut väärin vain yhden merkin. Punaisella merkityjä sanoja löytyi 9 kappaletta. Keltaisella merkityjä sanoja löytyi 17 kappaletta.

vatus näin toisiinsa. Meidän katsomuksemme mukaan on yksilön kaikinpuolinen kehittäminen samalla parasta yhteisökasvatusta. Tämän mukaan on kasvatuksemme lähimmäääränä verevän suomalaisen ihmisen ja kunnon kansalaisen kehittäminen. Tätä rää me pidämme koetinkivenämme kasvatustaljemme ja toimenpiteitemme oikeutusta ja arvoa määrättäessä."

Kuten edellisestä ilmenee, tarkastelee maisteri Kaljumen artikkelissaan kasvatuskysymystä humanistisen etiikan pohjalta. Ihmisen uskonnolliset tarpeet, joita kirjoituksessa ei erikseen mainita, saavat tällöin enintään tasa-arvoisina toisten tarpeiden rinnalla sijan kasvatustaljemmassa. Uskonnollisen elämän keskeinen merkitys persoonallisuuden muodostuksessa tulee siten kokonaan syrjäytetyksi (vrt. 1 Tim. 4: 8: "jumalasuudesta on hyötyä kaikkeen, koska sillä on elämän lupaus, sekä nykyisen että tulevaisen"), puhumattakaan siitä menetyksestä, mikä kasvatuksen päämääriä on kohdannut niiden menetettyä tuonpuoleiset ikuisuustavoitteensa. Vrt. myös sotilaspastori Raivion selostusta piispa Runestamin kirjasta "Kxisten-dom och världsforbättring" lehtemme tässä numerossa

Tässä yhteydessä ei voi olla viittaamatta siihen, (sana puuttuu) arvon jo pakanat ovat ymmärtäneet antaa uskonnolle kansan elämäsä ja sen kasvatustaljemmassa. Niinpä keisari Neron aikana elänyt Petronius Arb i ter kirjoittaa kirjassansa Trimalkion pidot: "Vlitäi tästä tuleekaan, jolleivät jumalat eivätkä ihmiset sääli tätä kauipunkia! Niin totta kuin tahdon saada iloa lapsistani, luulen, että kaikki (sana puuttuu) (nimittäin tavaton kuivuus) tulee jumalilta. Sillä kukaan ei pidä taivasta taivaana, ei kukaan noudata paastoa, ei kukaan välitä hi-iiuistakaan Juppiterista, vaan kaikki vain katsovat alaspäin ja laskevat rahojaan.

Ennen aikaan menivät naidut naiset pitkissä vaipoissaan paljain jaloin vuorelle, hiukset hajallaan, mielet puhtaina, ja rukoilivat Juppiterilta vettä. Ja niinpä satoikin heti kuin saavista kaataen. Joko silloin tai ei koskaan. Ja kaikki palasivat kotiin märkinä kuin koirat. Senfvuoksi on jumalilla jalat villaan kää-ittyinä (s.o. he pitivät yllä läimmintä ja kuivaa ilmaa), että me emime enää ole hurskaita." (Trimalkion -pidot, siv. 68-69.)

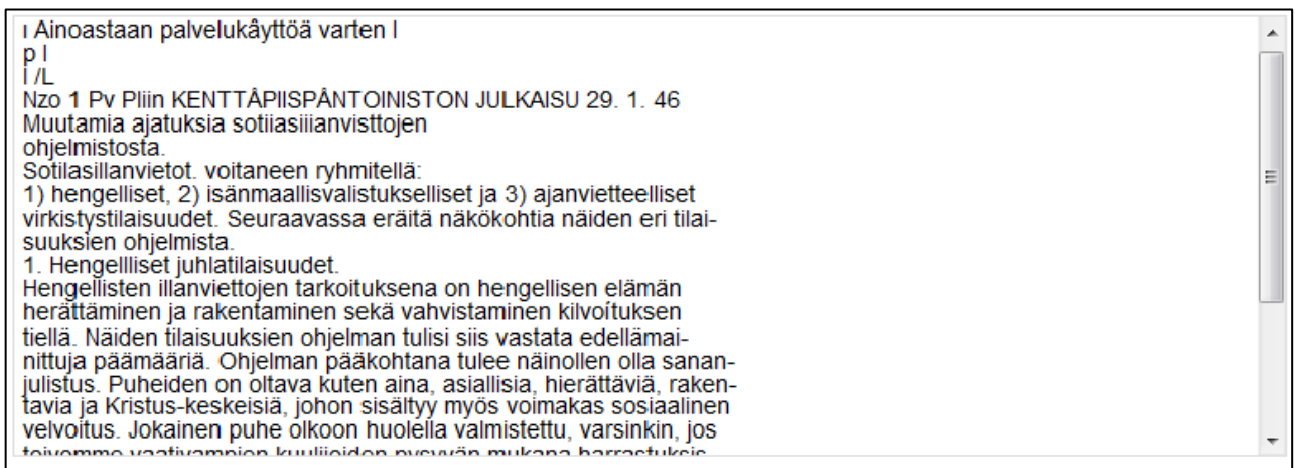
15

Kuva 9. i2OCR-ohjelman tekstintunnistusnäyte testiaineistosta 2

Ohjelma pystyy tulostamaan testatun tiedoston vain Word-tiedostoksi tai tekstiksi selaimen. Virheiden korjaaminen kestää kauan vertaamalla alkuperäis-tekstiin.

5.1.4 New OCR -ohjelma

New OCR -ohjelmassa kieleksi voi valita suomen, englannin ja ruotsin lisäksi 26 kieltä (Taulukko 1, s.16), mutta eräessä testissä vaikka ohjelmassa /5/ valittiin kieleksi suomi, ä ja ö eivät silti aina tulostuneet, vaan joissain kohdin ne olivat a ja o. Käsiteltävän tiedoston koko saa olla enintään 5 MB (Megabyte) eli megatavua. Tavu on tallennuskapasiteetin mittayksikkö. Testissä (Kuva 10) käytettiin testiaineistoja 1 ja 2 (Taulukko 2, s.17).



Kuva 10. New OCR -ohjelman tekstintunnistusnäyte

Seuraavanlaisesti tulostui käännös (Kuva 11) testiaineistosta 2 (Taulukko 2, s.17); Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat sanat, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisella merkittyjä sanoja löytyi 19 kappaletta. Keltaisella merkittyjä sanoja löytyi 9 kappaletta.

Sukupuolitautilien leviäminen ja taistelu niitä vastaan.

(Muistiinpanoja tri Eero Uroman Ilmavoimien pasitareiden neuvottelukokouksessa 25. 1. 46 pitämästä esitelmästä)

Sukupuolitautilikysymys on tällä hetkellä erittäin ajankohtainen ja maamme kannalta sargent tärkeä. Numerotietoja ei kuitenkaan ole riittävästi käytettävissä. V. 1938 oli maassam,m.e n. 900 uutta* kuppatautitapausta. Ne vähenivät jossakin määrin tultaessa vuoteen 1940. jolloin kansaim-me eli määrättyllä tavallatinoraalisesti korkeata aikaa. Silloin oli kuppatautitapauksia vain n. 800, gonorrheatapauksia 2000. Uuden sodan alkaessa sukupuolitaudit. alkoivat vefelhitellexi levitä. V. 1942 oli syfilistautitartuntoja 2.500. Lukumäärä mutsi nopeasti. V. 19-13, jolloin saimme uuden sukupuolitautilain, tautitapauksien lukumäärä oli noussut ä.000:ec-n. Kuitenkin rupesi kippatauti v. 1943 osoittamaan liar:il~:e;m.f.sta. V. 1944 oli tapauksia jälleen vain n. 3.000. Mutta v 1945 se saavutti jäll lc:~n uuden maksimln, 5.500.

Gnnorrhea-tapaukset ovat koko ajan lisääntyneet hitaasti. V. 1941 oli niitä enemmän kuin v. 1938, ts. n. 5 ä 6 tuhatta Mutta v. 1945 lukumäärä no-usi 22.000:een.

Aikaisemmin. hai.se^t mu-odustivat n. kolmanneksen sukupuolitautiltapauksista. V. 1943-45 ne ovat huomattavasti ylittäneet mies-tapausten luikumäärän. Mutta vzn 1945:n 22.000:sta gonorrh~ea-tapauksesta tuli vain 6.000 naisten osalle. VI ist~ähän tämä johtuu? Ei _17 -i

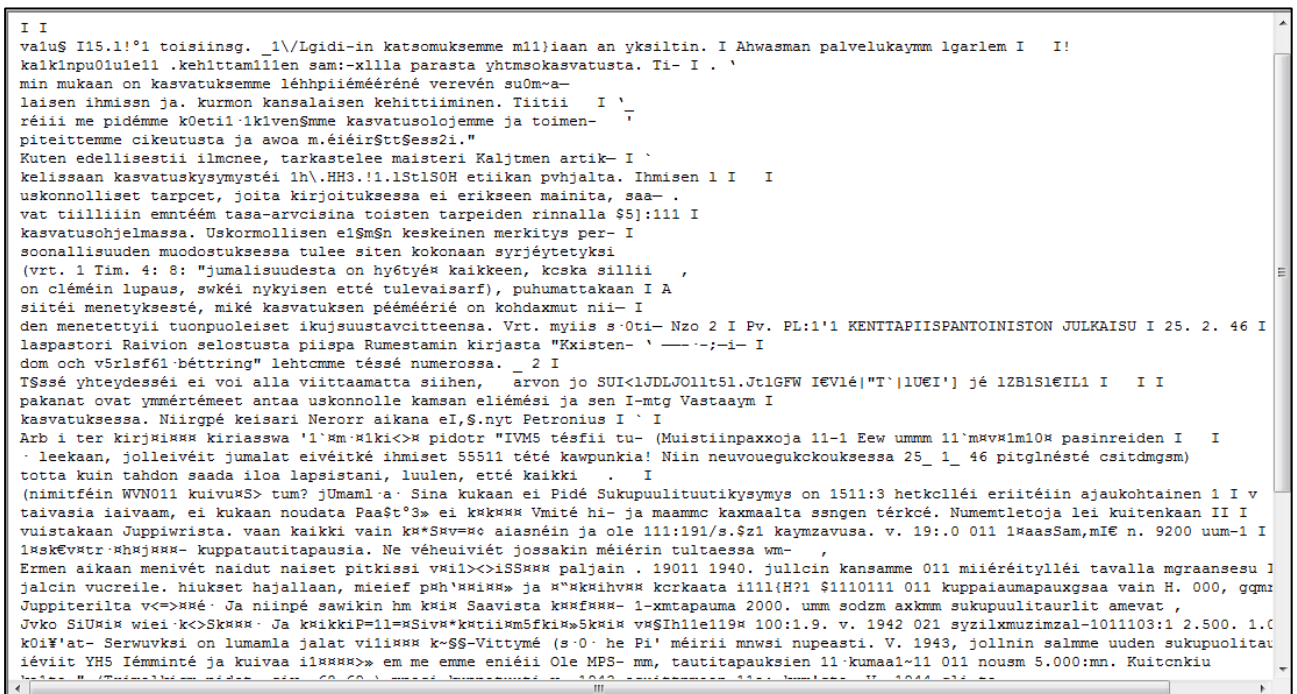
Kuva 11. New OCR -ohjelman tekstintunnistusnäyte testiaineistosta 2

Ohjelman tulostamassa Word-tiedostossa oli erittäin paljon virheitä eli tunnistamattomia merkkejä niin paljon, että virheiden korjaaminen olisi vienyt liian paljon aikaa.

5.1.5 OCR Convert -ohjelma

OCR Convert -ohjelman /6/ kielivaihtoehtoina löytyvät espanja, englanti, saksa, italia ja hollanti (Taulukko 1, s.16), eli suomen merkkejä ei tästä tuotteesta löydy.

Tulokset aukeavat tekstinä uuteen välilehteen. Ensimmäisessä testitiedostossa oli edes muutama sana joista sai selvää (kuva 12), mutta toisella kerralla tuloksena oli vain muutama rivi, joista ei saa selvää.



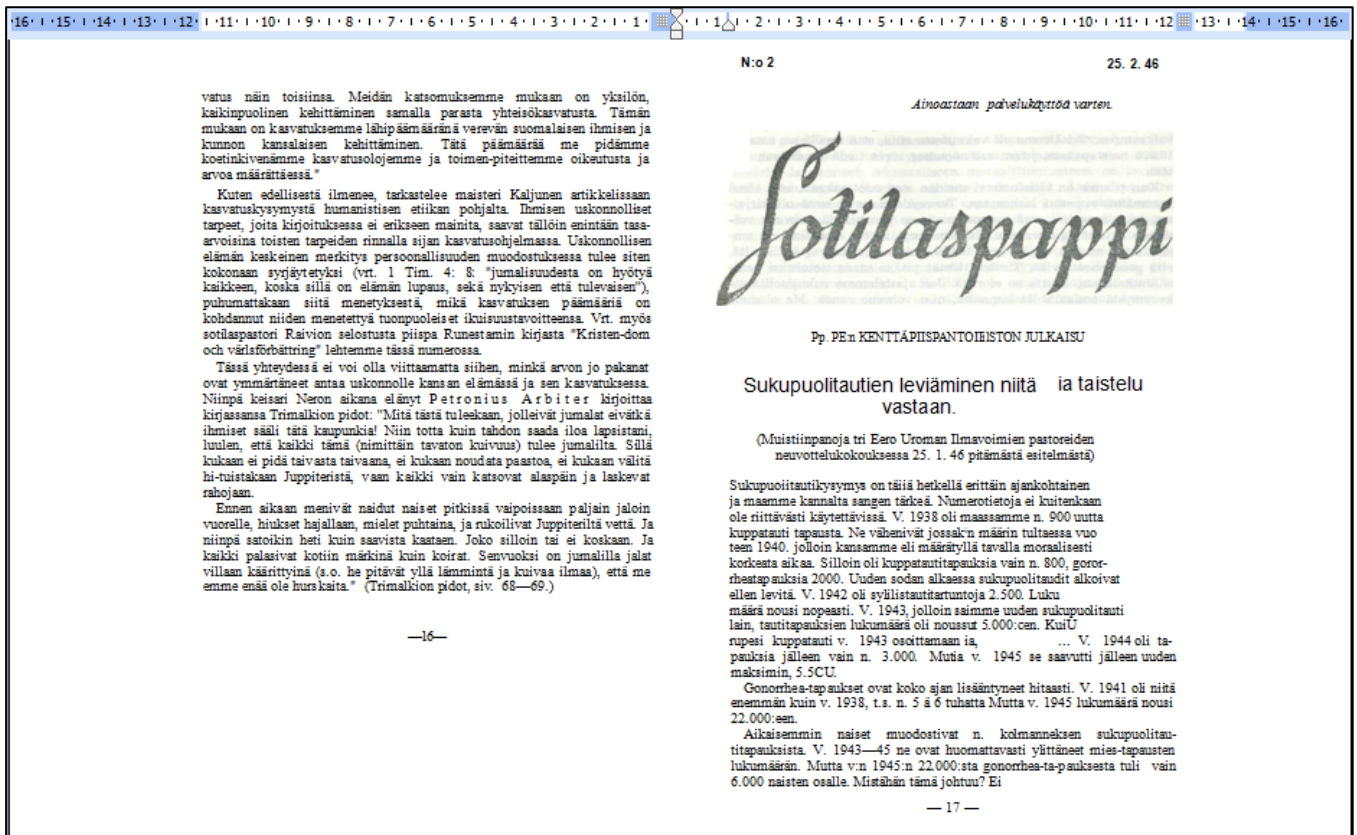
Kuva 12. OCR Convert -ohjelman tekstintunnistusräjä

OCR Convert -ohjelman testissä käytettiin testiaineistoa 2 (Taulukko 2, s.17). Tulos oli käyttökelpoisuudeltaan heikko. Lähes jokainen sana oli tulostunut väärin, eikä tämän työn tarkoituksena ole kirjoittaa koko tekstiä uudelleen.

5.1.6 OCR Now! -ohjelma

OCR Now! -sivusto vaatii rekisteröitymisen /7/, jotta kuvia pääsee käsittelemään. Kieliä löytyy englannin ja suomen lisäksi 14.

Lopputulokset Wordilla olivat tyydyttävät (Kuva 13). Monia virheitä silti löytyi. Tiedostojen käsittely vaatii krediittejä, joita alussa on vain kaksi, joten kovin ilmaiseksi tätä sivustoa ei voi kutsua. Testissä käytettiin testiaineistoa 2 (Taulukko 2, s.17).

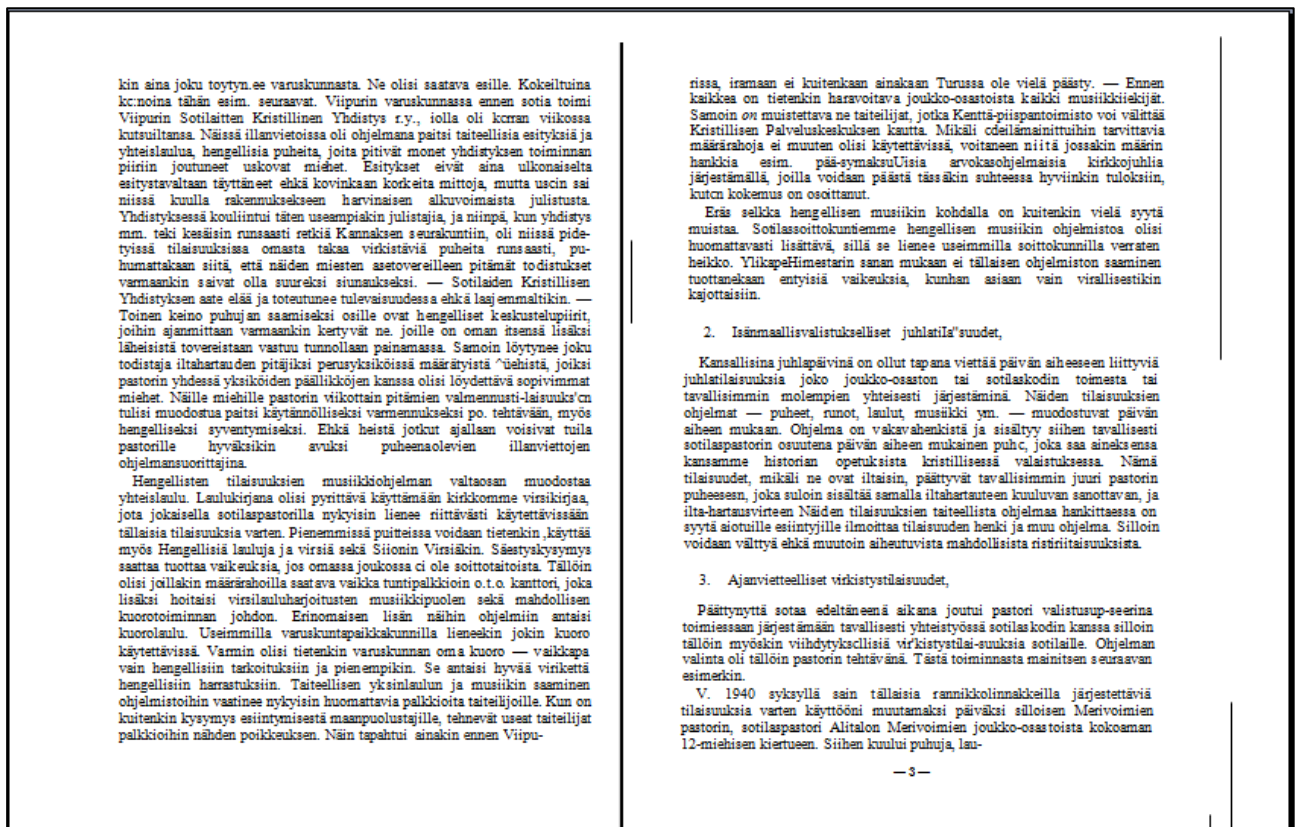


Kuva 13. OCR Now! -ohjelman tekstintunnistusnäyte

Krediitit maksavat (Taulukko 5, s.18) 20 kpl 1.99 puntaa ja 500 krediittiä (Credit) taas on 31,99 puntaa. Krediitit ovat Internet-sivustoilla käytettävää valuuttaa. Virheiden korjaus veisi niiden pienen määrän huomioon ottaen vain vähän aikaa.

5.1.7 OCR online -ohjelma

OCR online -ohjelma vaatii rekisteröitymisen /8/. Ilmaiseksi kuvatiedostoja voi käsitellä vain 10 kappaletta viikossa, jonka jälkeen ohjelma muuttuu maksulliseksi. Kielivaihtoehtoina (Taulukko 1, s.16) listassa ovat englanti, tšekki, ranska, saksa, italia, latina, latvia, puola, portugali, romania, venäjä, espanja ja turkki. Lisäksi OCR Onlinen kielivaihtoehtona on mikä tahansa kieli. Ohjelmassa on asetuksista lisättävä suomen kieli, koska listassa on aluksi vain 15 kieltä. Yhteensä kieliä on 153. Vilkaisulla ei käännöksestä (kuva 14) virheitä paljoa huomaa, mutta tarkemmin katsoessa virheitä näkyy useita.



Kuva 14. OCR Online -ohjelman tekstintunnistusnäyte

Ohjelmalla on mahdollisuus muuttaa kuvatiedosto (Taulukko 3, s.17) neljään (4) eri tiedostomuotoon (Taulukko 4, s.18). Testissä käytettiin testiaineistoa 3 (Taulukko 2, s.17). Virheiden korjaaminen niiden pieni määrä huomioon ottaen käy nopeasti.

5.1.8 OCR terminal -ohjelma

OCR terminal -ohjelma liittyi myöhemmin ABBYY Finereader Onlineen /9/.

Kieliä tässä ohjelmassa on suomen lisäksi pari kymmentä (Taulukko 1, s.16). Ohjelma vaatii rekisteröitymisen. Rajoituksena on 20 ilmaista sivua kuukaudessa (Taulukko 5, s.18). Ohjelma tukee TIFF-, PDF-, JPEG- ja muita kuvaformaatteja (Taulukko 3, s.17). Ohjelma muuttaa kuvan neljään (4) eri tiedostomuotoon (Taulukko 4, s.18).

Seuraavanlaisesti tulostui käännös (Kuva 15) testiaineistosta 4 (Taulukko 2, s.17); Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat virheet, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisia virheitä löytyi 1 kappaletta. Keltaisia virheitä löytyi 1 kappaletta.

keskinäiseen polemiikkiin, vaan sen sijaan olisi pyrittävä siihen hedelmälliseen yhteistyöhön, jota usko ja psykoterapeuttinen tekniikka sielunlääkinnälle tarjoavat. Niinpä voi uskonnollinen kokemus äkkiä avata potilaan silmät näkemään synnin synniksi, mikä puolestaan auttaa häntä löytämään sen lähteenä olevan kompleksin, ja toisaalta saattaa jälleen psykoanalyttinen käsittely selvittää hänelle hänen kompleksinsa, mikä voi johdattaa hänet tuntemaan siitä aiheutuvan synnin, jolloin hänelle avautuu myös Jumalan anteeksiantamuksen tie. Lääkärin on vain tehtävä oikea diagnoosi voidakseen valita sen edellyttämän joko hengellisen tai psykoterapeuttisen hoitotavan.

Tekniikan ja uskon välinen vuorovaikutus tapahtuu tällöin seuraavasti. Tekniikka, psykoanalyysi, tutkii ja tuo päivänvaloon ihmissielun probleemit, usko, armo, hävittää ne tavalla, jota emme osaa edes selittää. Tekniikkaa voidaan oppia, ja se on käyttökelpoinen maaperän muokkauksessa, mutta hengellisiä kokemuksia sillä ei kyetä aikaansaamaan. Kun Tournier'ltä on pyydetty sielunhoito-opillisia neuvoja, hänen on täytynyt tunnustaa, ettei hän osaa niitä antaa, sillä hän ei tiedä, kuten hän sanoo, mitä on tehtävä todellisessa sielunhädässä olevan ihmisen auttamiseksi. Hän on suorastaan sitä **<m:eltä**, ettei tuollaisessa tapauksessa ole tehtävä mitään, sillä jos ihminen siinä johonkin ryhtyy, se jo on itseensä eikä Jumalaan **luottamsta**. Siksi siinä ei olekaan tähdellistä se, mitä me teemme, vaan mitä me olemme. Jos itse olemme yhteydessä Jumalaan, ja tämä yhteys opettaa meidät näkemään viheliäisyytemme ja oman kykenemättömyytemme auttaa toisia, silloin jo tällä paikalla olemisemme on omiaan johtamaan heidät uskonnollisiin kokemuksiin. Tekniikkaan nähden voimme toisia neuvoa, mutta kun usko on kysymyksessä, ei neuvoistamme ole mitään hyötyä. Tekniikka auttaa ihmistä kuitenkin tutkimaan itseään ja syventämään hänen moraalisia arviointiperusteitaan. Tällöin ei normina ole enää muodollinen moraalitapa, joka pitää silmällä vain sitä, mitä ihminen tekee, vaan syvällisempi moraalikäsitelmä, joka kysyy, minkälainen ihminen on. Tässäkin voi psykoanalyysi auttaa meitä löytämään kanssaihmisistämme heidän sydämensä usein salassa olevat hyvät ominaisuudet, ja mitä tulee meihin itseemme, se on omiaan paljastamaan sisimpämme pohjamutiaa myöten.

Muodollinen moraalitapa, jonka normina siis on, onko jokin teko itsessään hyvä vai paha, onko se hyödyllinen vai ei ja mitä ihmiset siitä sanovat, johdattaa orjuuttavaan ja seurauksiltaan monta kertaa tuhoisaan moralismiin. Elämällä on kuitenkin lakinsa, ja siksi esim. kirkon on asetettava moraalille tietyt vaatimukset, sillä kuinka hyvät ihmisen tarkoitusperät ovatkin, paha vaanii häntä yhtenä ja saa hänet helposti saaliikseen, ellei kirkko ole opetuksellaan häntä opastamassa. Mutta tämän moraalin normina tulee olla Jumalan

— 50 —

Kuva 15. OCR Terminal –ohjelman tekstintunnistusnäyte testiaineistosta 4

Tuloksista löytyi erittäin vähän virheitä Wordissa. Jotkut kirjaimet ovat muuttuneet eikä teksti ole samalla tasolla kuin kuvassa, joten väliviivoja voi olla alkuperäisen kuvan takia joissain kohdissa, missä niitä ei pitäisi olla. Kuvan kuitenkin täytyy olla riittävän iso, jotta ohjelma tunnistaa merkit. PDF-tiedostossa ei sanahauulla tunnista hajanaisia sanoja.

5.2 Koneelle asennettavia OCR-ohjelmia

Kun testattavat selainpohjaiset ohjelmat oli käyty läpi, testattiin tekstintunnistusohjelmia, jotka asennetaan työasemalle. Testattavia ohjelmia kokeiltiin viisi (5) kappaletta ja tutkittiin, kuinka hyvä lopputulos niillä saadaan.

Testiaineistoista (Taulukko 7) pääasiassa testiaineistoa 2 käytettiin, mutta vaihtelevuuden vuoksi käytettiin myös testiaineistoja 1 ja 5. Useimmat testattavista ohjelmista olivat vain kokeiluversioita, joiden käyttöikä oli rajallinen, ja joissakin ohjelmissa myös tallennettavien tiedostojen lukumäärä oli rajattu. Yhteen työhön keskitettynä rajallinen aika kuitenkin riittää. Lukumääräisesti rajattujen ohjelmien ongelma on, että tulostettavia tiedostoja voi olla liikaa. Ohjelmissa oli eroavaisuuksia tuettujen kielten kanssa (Taulukko 6). Jotkut ohjelmat tukivat monia kieliä ja jotkut tunnistivat vain englanninkielestä löytyviä merkkejä ja sanoja. Ohjelmat erosivat myös siinä, mitä kuvaformaattia ne tukevat (Taulukko 8) ja mihin tiedostomuotoon käsitelty tiedosto tulostuu (Taulukko 9). Tuettuja kuvaformaatteja oli runsaasti: PDF, GIF, TIFF, JPEG ja PNG, BMP ja PPT (Power Point Table), joka on Microsoftin PowerPoint-ohjelman käyttämä tiedostomuoto ja PCX (Personal Computer eXchange), joka on vähenevässä käytössä oleva kuvaformaatti.

Tiedostomuotoja, jotka voi valita käsittelyn jälkeiseksi formaatiksi, oli myös muutamia: DOC, PDF, TXT, RTF, XLS ja HTML (Hypertext Markup Language), joka on hypertekstin merkintäkieli.

Ohjelma	suomi	englanti	ruotsi	muut kielet yht.
Acrobat pro 10	x	x	x	25
Autobahn		x		
Free OCR 3.0		x		
Finereader 10	x	x	x	186
Smart OCR		x		
PDF Converter 7		x		

Taulukko 6. Työasemalle asennettavien tekstintunnistusohjelmien kielet

Käsittelyissä käytettiin kolmea (3) eri testiaineistoa (Taulukko 7), jotta saataisiin erilaisia näkökulmia ja vaihtelevuutta.

Ohjelma	Testiaineisto 1	Testiaineisto 2	Testiaineisto 5
Acrobat pro 10	x		
Autobahn		x	
Free OCR 3.0		x	
Finereader 10		x	
Smart OCR		x	
PDF Converter 7			x

Taulukko 7. Ohjelmissa käytetyt testiaineistot

Ohjelma	PDF	GIF	TIFF	JPEG	BMP	PNG	PPT	PCX
Acrobat pro 10			x	x				
Autobahn	x		x	x				
Free OCR 3.0	x	x	x	x	x			
Finereader 10	x	x	x	x	x	x		x
Smart OCR	x	x	x	x	x	x		
PDF Converter 7	x	x	x	x	x	x	x	

Taulukko 8. Työasemalle asennettavien tekstintunnistusohjelmien käsiteltävät kuvaformaattit

Ohjelma	DOC	PDF	TXT	RTF	XLS	HTML
Acrobat pro 10		x				
Autobahn		x				
Free OCR 3.0	x		x			
Finereader 10	x	x	x	x	x	
Smart OCR	x	x	x	x		x
PDF Converter 7		x				

Taulukko 9. Työasemalle asennettavien tekstintunnistusohjelmien tulostettavat tiedostomuodot

Ohjelma	Ei rajoituksia	Hinta (€)	Kokeilu-aika (vrk)	Tiedostorajoitus (kpl)	Vesileima
Acrobat pro 10		199 dollaria	30		
Autobahn		960 euroa			x
Free OCR 3.0	x				
Finereader 10		129 euroa	15	50	
Smart OCR		99,90			x
PDF Converter 7		99,99 dollaria	30		

Taulukko 10. Työasemalle asennettavien tekstintunnistusohjelmien rajoitukset

Rajoituksettomia ohjelmia oli vain yksi, mutta ohjelman tuottaman tunnistuksen laatu ei ole onnistunut. Paras lopputulos oli Finereader 10:llä, mutta tiedostorajoitusten takia projektia ei saanut ohjelmalla tehtyä. Toiseksi paras ohjelma oli PDF Converter 7, jonka kuukauden koeaika oli riittävä selvityksen tekoon.

5.2.5 Abby Finereader 10 -ohjelma

Abby Finereader 10 -ohjelmasta löytyy 189 kieltä (Taulukko 6), joka on paljon. Ohjelman /10/ tunnistuksesta löytyi erittäin vähän virheitä.

Ohjelman ilmaisversiolla (Taulukko 10) ei saa kuin 50 sivua tallennettua, eikä se toimi kuin 15 vuorokautta. Ohjelma olisi vaihtoehto työn suorittamiselle, mutta rajoitustensa takia ohjelmaa ei voi käyttää. Aivan virheettömästi ohjelma ei tekstiä tunnistanut, mutta virheitä oli erittäin vähän. Ohjelman hinta on 129 euroa.

Jos joku haluaisi kyseisellä ohjelmalla vastaavan työn suorittaa ja haluaisi yhdistää tiedostot yhdeksi, PDF-tiedostot saa helposti yhdistettyä selaimesta toimivalla ilmaisella Merge PDF -ohjelmalla (Kuva 16) /11/. Myös Word-tiedosto oli riittävän onnistunut.

The screenshot shows the Merge PDF website interface. At the top left, there is a heading "Merge PDF Documents Easily and for Free !" followed by a paragraph explaining the service. Below this, there are several links: "Ads by Google", "Merge PDF", "Free PDF", "Word to PDF", "Convert PDF", and "PDF Maker". On the right side, there is a section for "AVS Document Converter" listing various file formats: EPUB, PDF, DOC, DOCX, ODT, FB2, DjVu, XPS, HTML, MHT, PPT, PPTX, RTF, TXT, TIFF, GIF, JPEG, and PNG. A "Download Now" button with a downward arrow is present, along with the website URL "www.av4you.com".

Combine PDF Documents:

STEP 1
Select two (or more) PDF files to be joined together by pressing "Browse"

BROWSE

STEP 2
Press the "Merge PDF" button to upload and combine the selected documents!

MERGE PDF

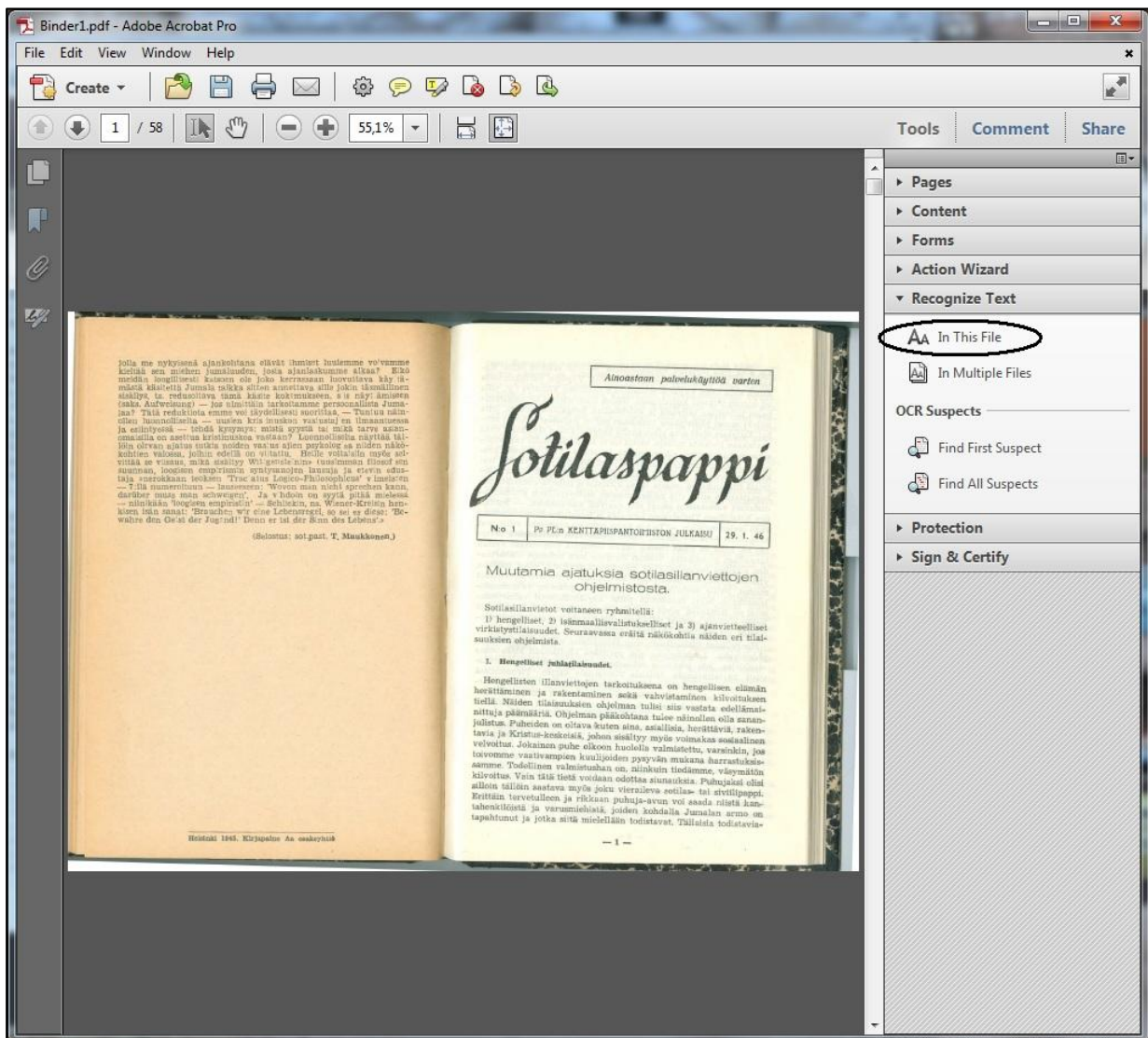
+1 527 Tweet 109
Recommend 2,290 people recommend this. Be the first of your friends.

Kuva 16. Merge PDF -selainpohjaisen ohjelman näkymä

Ohjelman käsittelemien tiedostojen lopputulokset olivat selvityksen ohjelmista parhaat. Virheitä oli erittäin vähän eikä niiden korjaamiseen mene kauaa.

5.2.6 Adobe Acrobat pro 10 -ohjelma

Adobe Acrobat pro 10 -ohjelmasta oli käytössä kolmenkymmenen (30) päivän kokeiluversio /12/, mutta toisin kuin Finereaderissa (k. 5.2.1), tässä ohjelmassa ei ollut rajoituksia dokumenttien tallennusten suhteen (Taulukko 10 s.33). Ohjelman alussa muutettiin ensimmäiseksi haluttu määrä skannattuja kuvia eli tässä tapauksessa 46. vuosikerta yhdeksi PDF-tiedostoksi, jonka jälkeen käsitettiin ohjelmaa tunnistamaan teksti (kuva 17).



Kuva 17. Adobe Acrobat pro 10 -ohjelman tekstintunnistusnäyte

Tekstintunnistuksen valmistuttua tarvitsi vain tallentaa tiedosto halutun nimiseksi ja vuosikerta oli valmis. Testissä käytettiin testiaineistoa 1 (Taulukko 7, s.32). Virheitä oli useita, eikä ohjelma löytänyt vuosikerran lukuja. Ohjelma kääntää kuvan vain PDF-tiedostoksi, joten virheitä ei voi korjata.

5.2.7 Autobahn DX -ohjelma

Autobahn DX -ohjelmalla /13/ pystyy yhdistämään niin PDF-, kuin TIFF-tiedostoja. Ennen tekstintunnistusta tuli muistaa valita OCR options -valikosta Searchable PDF -vaihtoehto. Sen valitsemalla ohjelma (Kuva 20) ymmärtää tunnistaa tekstin kuvan sisällä.

Lopputulos ei ollut onnistunut, ja vaikka olisi ollutkin, Wordissa (Kuva 18) tekstissä oli paljon virheitä ja PDF:n (Kuva 19) päällä oli vesileima, joka johtuu ohjelman kokeiluversiosta (Taulukko 10, s.33). Vesileima ei silti vaikuttanut päätökseen olla käyttämättä ohjelmaa selvityksessä, sillä lopputulos ei ollut riittävän hyvä.

Sukupuolitautiin leviämisen ja taistelu niitä vastaan.

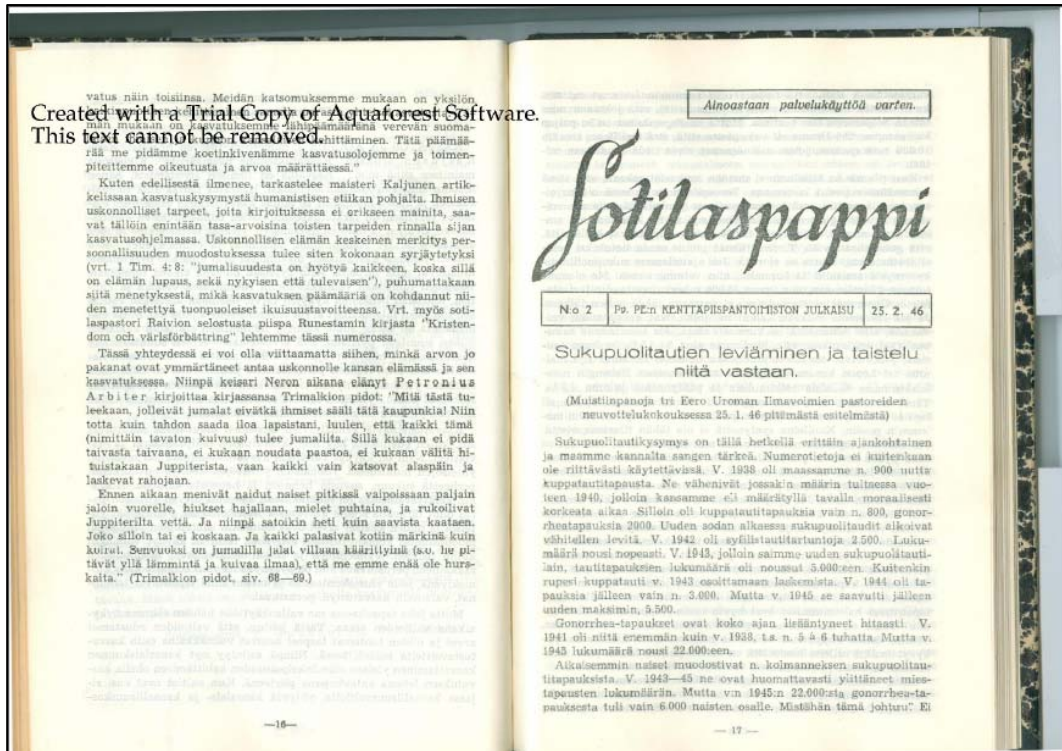
(Muistinpöytätri Eero Uroman Ilmavoimien pastoreiden
neuvottelukokouksessa 25. 1. 46 pitämästä esitelmästä)

Sukupuolitautiin leviäminen on tällä hetkellä erittäin ajankohtainen ja maamme kannalta sängin tärkeä. Numerotietoja ei kuitenkaan ole riittävästi käytettävissä. V. 1938 oli maassamme n. 900 uutta leikkauksetapauksia. Ne vähenivät jossakin määrin tultaessa vuoteen 1940, jolloin kansamme ei enää käytetty tavalla, moraalisesti korkeasta ajasta. Silloin oli leikkauksetapauksia vain n. 800, gonorrhoeatapauksia 2000. Uuden sodan aikana leikkauksetapaukset alkoivat vähitellen lisääntyä. V. 1942 oli syödiä tautitapauksia 2.500. Luku nousi nopeasti V. 1943, jolloin saimme uuden sukupuolitautiin taisteluun, tautitapauksien lukumäärä oli noussut 5.000:een. Kuitenkin rupe i leikkaukset v. 1943 osoittamaan laskeutumista. V. 1944 oli leikkauksetapauksia jälleen vain n. 3.000. Mutta v. 1945 se saavutti jälleen uuden maksiminsa, n. 5.500.

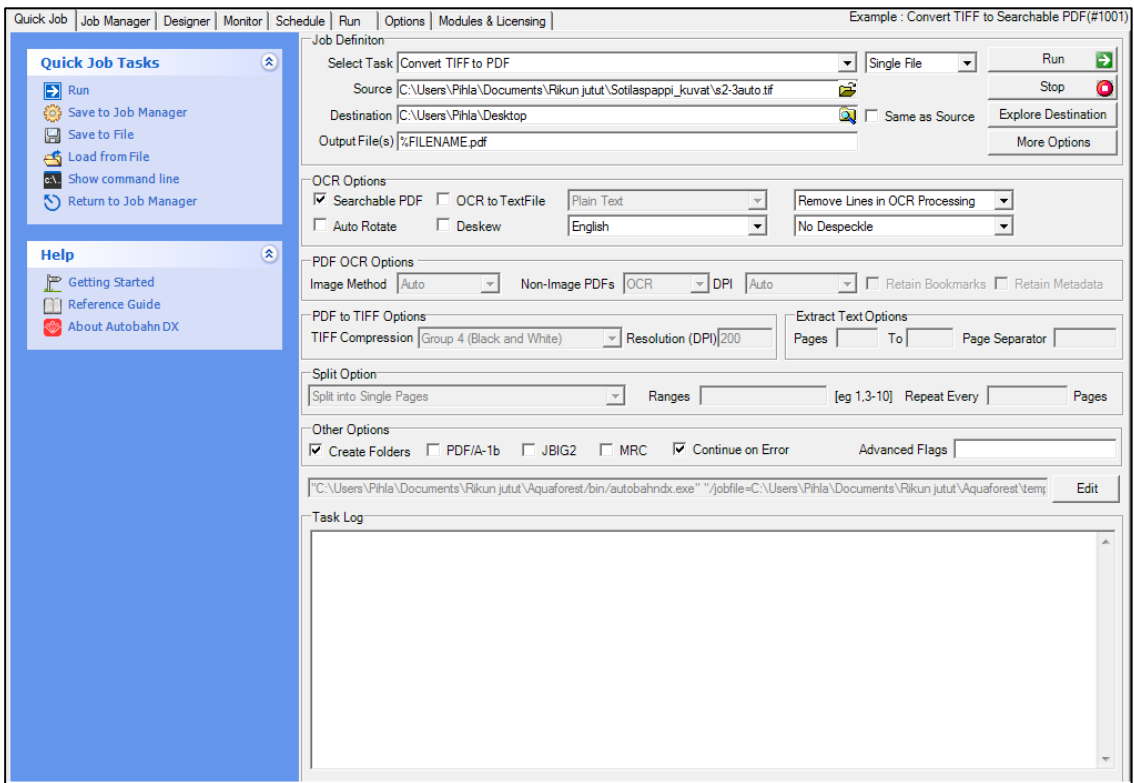
Gonorrhoea-tapaukset, ovat koko ajan lisääntyneet hitaasti V.
1941 oli niitä enemminkin kuin v. 1938, t. a. n. 506 tuhatta. Mutta v.
1945 luku nousi 22.000:een.

Aikaisemmin naiset muodostivat n. kolmannen sukupuolitautiin tapauksista V. 1943 — 45 ne ovat huomattavasti ylittäneet miesten tapauksien lukumäärän. Mutta v. n. 1945 n. 22.000:sta gonorrhoea-tapauksesta tuli vain 6.000 naisten osalle. Miesten osalle johtuu. Ei

Kuva 18. Autobahn DX:n .doc-version tekstintunnistusnäyte



Kuva 19. Autobahn DX:n PDF -version näkymä

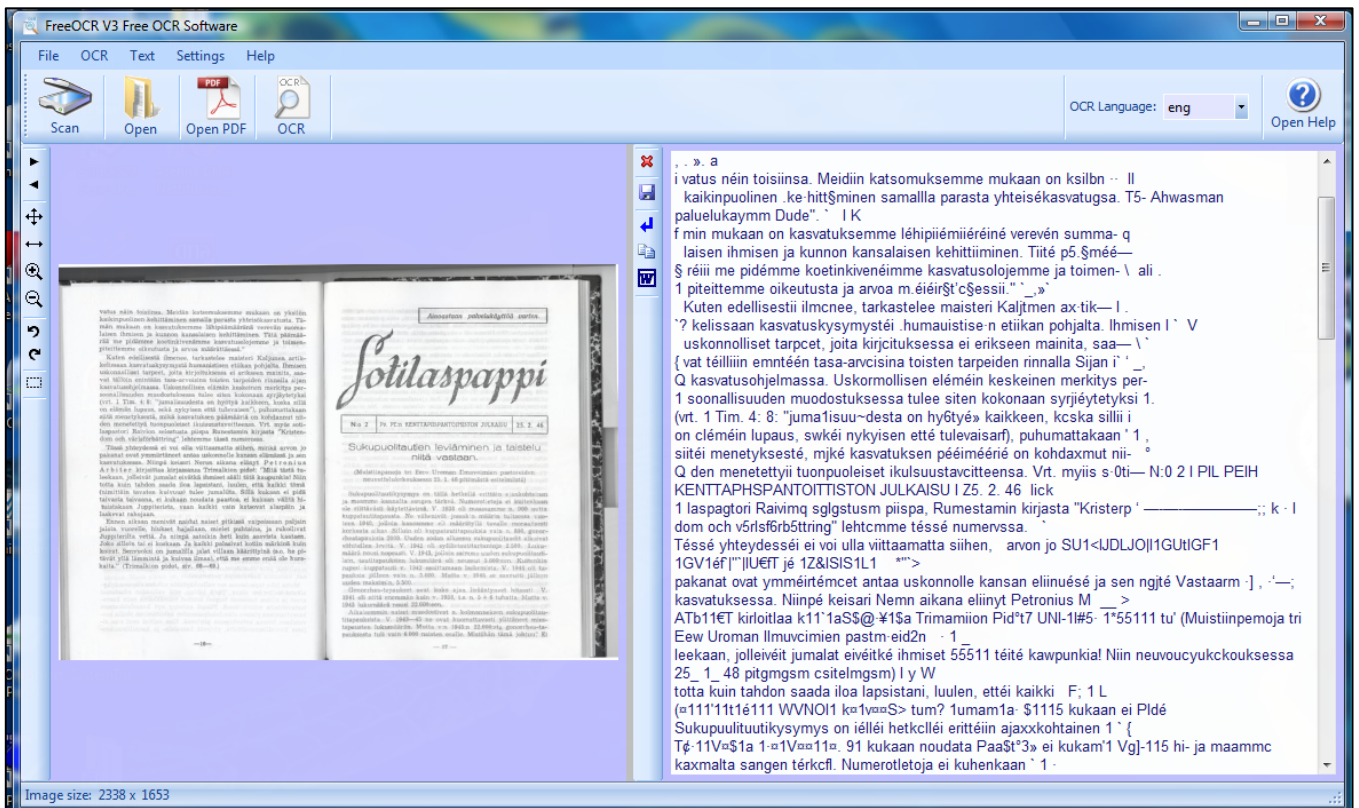


Kuva 20. Autobahn DX -ohjelman näkymä

Testissä käytettiin testiaineistoa 2 (Taulukko 7, s.32). Teksti vaati paljon korjauksia, mikä olisi vienyt liikaa aikaa.

5.2.8 Free OCR 3.0 -ohjelma

Free OCR 3.0 -ohjelma /14/ tuki vain englannin kieltä, joten se ei tunnistanut ääjiä ja ö-kirjaimia (Kuva 21), mutta tunnistustaso ei ole ohjelmassa muutenkaan onnistunut, joten tämä ongelma ei ole tärkeä.



Kuva 21. Free OCR 3.0 -ohjelman tekstintunnistusnäyte

Yhdestä aukeamasta ohjelma tunnisti muutaman sanan. Testissä käytettiin testiaineistoa 2 (Taulukko 6, s.31). Ohjelman laatu ei tunnituksen perusteella ollut lainkaan riittävä. Tekstin korjaamiseen olisi käytännössä mennyt yhtä kauan kuin koko tekstin uudelleen kirjoittamiseen.

5.2.9 Smart OCR -ohjelma

Smart OCR -ohjelma /15/ siirsi kuvan haluttuun tiedostomuotoon, mutta ei tunnistanut merkkejä riittävästi ja sivustolla lukee selvästi, että ohjelma jättää vesileiman (Taulukko 10, s.33) lopputulokseen (Kuva 22), joten pelkästään sen takia ohjelma ei täytä vaatimuksia.



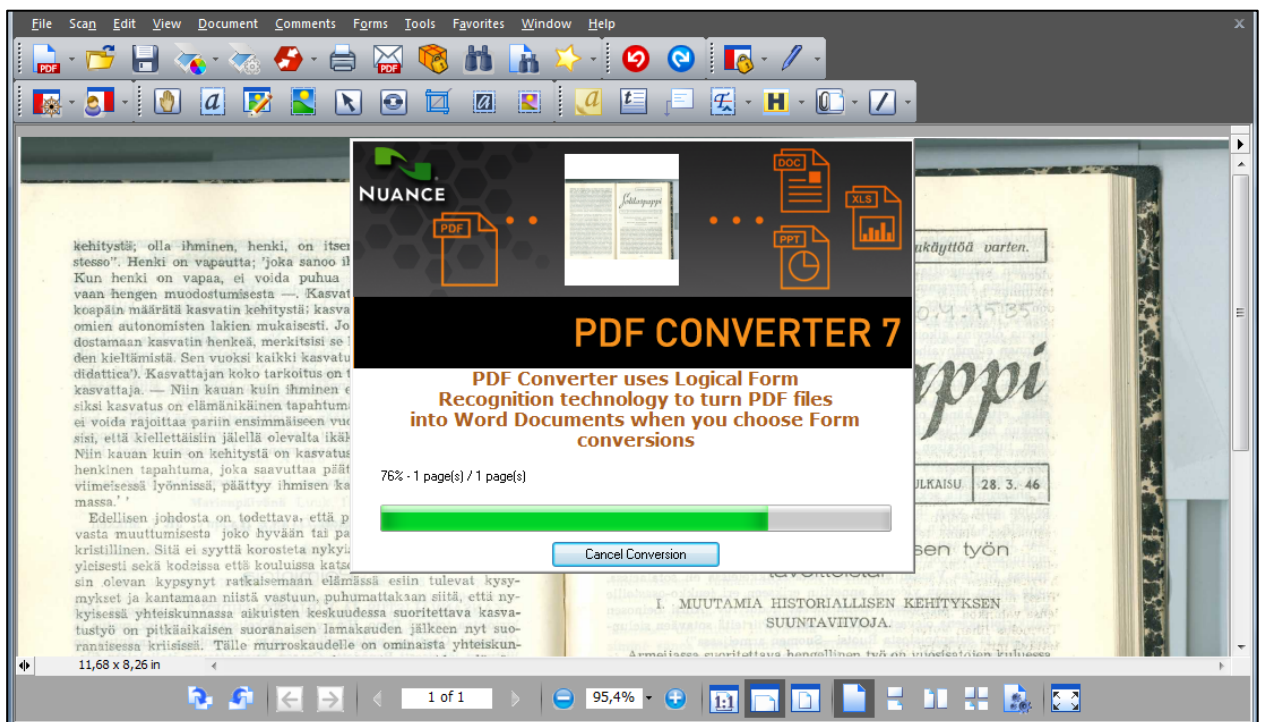
Kuva 22. Smart OCR -ohjelman tekstintunnistusnäyte

Testissä käytettiin testiaineistoa 2 (Taulukko 7, s.32). Tunnistuksen suurin ongelma on se, että ohjelma ei edes tunnista koko tekstiä, vaan parhaimmillaan hieman yli puolet tekstistä tulostuu.

5.2.10 PDF Converter Professional 7 -ohjelma

PDF Converter Professional 7 -ohjelmasta oli tarjolla 30 vuorokauden kokeiluversio /16/.

Lopputulos Word-dokumentiksi muuttamisen jälkeen oli tällä ohjelmalla vaihtelevaa, sillä sanat, joita ohjelma ei Word-dokumentiksi muutettaessa täysin tunnistanut, tunnistettiin PDF-muotoon muutettaessa (Kuva 23). Ohjelman tuottama Word-dokumentti ei tunnistanut ä- eikä ö- kirjaimia, mikä olisi tietenkin suotavaa. Tämä on huono asia, sillä toisin kuin PDF-tiedosto, Word-dokumentti tunnisti harvalla kirjoitetut sanat. PDF-formaatissa ä- ja ö -kirjaimet eivät onneksi ole ongelma, koska PDF tunnistaa ä- ja ö- kirjaimet. Testissä käytettiin testiaineistoa 5 (Taulukko 6, s.31).



Kuva 23. PDF Converter 7 -ohjelman tekstintunnistus

Ohjelma pystyy yhdistämään kuvatiedostot yhdeksi PDF-tiedostoksi, mikä on suositeltavaa ennen tekstintunnistuksen suorittamista. Pakollista tämä ei ollut, mutta nopeutti ja helpotti tunnistusprosessia huomattavasti. Word-tiedostoksi kääntäminen onnistui vasta, kun kuvatiedostot olivat ensin muutettu PDF-muotoon. Ohjelma oli laadultaan selvityksen käytettävien, mutta silti puutteelli-

nen. Ensimmäisessä tarkastelussa virheitä ei huomattu monta, mutta tarkemman tarkastelun jälkeen virheitä löytyi useita.

6 Tulosten esittely

Opinnäytetyön teko aloitettiin skannaamalla Sotilaspappi-aikakausilehden vuosikerta aukeama kerrallaan. Seuraavaksi mietittiin sopiva testiaineisto, jossa olisi tarpeeksi erilaista tekstiä ja erilaisia merkkejä. Vaihtelevuuden ja kattavuuden vuoksi päätestiaineiston lisäksi valittiin vielä neljä testiaineistoa.

Ohjelmat etsittiin Internetistä ja kaikki löydetty ohjelmat testattiin, jonka jälkeen selvitettiin paras ohjelma eri ohjelmien virheiden määriä vertaamalla. Paras ohjelma oli sellainen, jossa oli vertailun jälkeen pienin määrä virheitä. Käyttökelpoisin vaihtoehto oli PDF Converter professional 7 -ohjelma, joka asennettiin työasemalle. Ohjelmalla tehtiin vuosikerran kuvatiedostoja yhdistämällä yksi PDF, jonka jälkeen valittiin ohjelmasta tekstintunnistus, joka tunnisti PDF:stä merkit. Tämän jälkeen tallennettiin vuosikerta-PDF ja nimettiin se (Sotilaspappi vuosikerta 46).

Testimateriaalina toimiva aikakausilehden vuosikerta oli käsin ladottua tekstiä, joka aiheutti sanoissa epätasaisuutta. Toinen vaikeuksia tuottava ongelma oli harvaan kirjoitetut sanat, joiden tunnistaminen osoittautui ongelmalliseksi. Kyseessä oli PDF-tiedosto, joten virheiden korjaaminen oli mahdotonta.

Vertailin tuloksia kolmen eri ohjelman välillä. Testiaineistona käytettiin testiaineisto 2:n (Taulukko 7, s. 32) ensimmäistä sivua. Vertailun ohjelmat valittiin niistä ohjelmista, jotka kääntävät lähdeaineiston suoraan PDF -tiedostoksi. Vertailussa olivat selvityksen parhaiten toimivaksi ohjelmaksi osoittautunut Abby Finereader 10 -ohjelma, selvityksen käyttökelpoisimmaksi ohjelmaksi osoittautunut PDF Converter Professional 7 -ohjelma ja Adobe Acrobat Pro 10 -ohjelma. Virheiksi ei merkitty ä- ja ö-pisteiden puuttumista, koska ne eivät estäneet sanojen löytymistä haussa.

Seuraavanlaisesti tulostui käännös (Kuva 24) Finereader 10 -ohjelmalla (k. 5.2.5, s. 34), Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat

virheet, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisia virheitä löytyi 0 kappaletta. Keltaisia virheitä löytyi 3 kappaletta.

vatus näin toisiinsa. Meidän katsomuksemme mukaan on yksilön, kaikinpuolinen kehittäminen samalla parasta yhteisökasvatusta. Tämän mukaan on kasvatuksemme lähipäämääränä verevän suomalaisen ihmisen ja kunnan kansalaisen kehittäminen. Tätä päämäärää me pidämme koetinkivenämme kasvatusohjelmamme ja toimenpiteittemme oikeutusta ja arvoa määrätessä."

Kuten edellisestä ilmenee, tarkastelee maisteri Kaljunen artikkelissaan kasvatuskysymystä humanistisen etiikan pohjalta. Ihmisen uskonnolliset tarpeet, joita kirjoituksessa ei erikseen mainita, saavat tällöin enintään tasa-arvoisina toisten tarpeiden rinnalla sijan kasvatusohjelmassa. Uskonnollisen elämän keskeinen merkitys persoonallisuuden muodostuksessa tulee siten kokonaan syrjäytetyksi (vrt. 1 Tim. 4: 8: "jumalisuudesta on hyötyä kaikkeen, koska sillä on elämän lupaus, sekä nykyisen että tulevaisen"), puhumattakaan siitä menetyksestä, mikä kasvatuksen päämääriä on kohdannut niiden menetettyä tuonpuoleiset ikuisuustavoitteensa. Vrt. myös sotilaspastori Raivion selostusta piispa Runestamin kirjasta "Kristendom och världsförbättring" lehtemme tässä numerossa.

Tässä yhteydessä ei voi olla viittaamatta siihen, minkä arvon jo pakanat ovat ymmärtäneet antaa uskonnolle kansan elämässä ja sen kasvatuksessa. Niinpä keisari Neron aikana elänyt **Petronius Arbitr** kirjoittaa kirjassansa Trimalkion pidot: "Mitä tästä tuleekaan, jolleivät jumalat eivätkä ihmiset sääli tätä kaupunkia! Niin totta kuin tahdon saada iloa lapsistani, luulen, että kaikki tämä (nimittäin tavaton kuivuus) tulee jumalilta. Sillä kukaan ei pidä taivasta taivaana, ei kukaan noudata paastoa, ei kukaan välitä **hutuistakaan** Juppiteristä, vaan kaikki vain katsovat alaspäin ja laskevat rahojaan.

Ennen aikaan menivät naidut naiset pitkissä vaipoissaan paljain jaloin vuorelle, hiukset hajallaan, mielet puhtaina, ja rukoilivat Juppiteriltä vettä. Ja niinpä satoikin heti kuin saavista kaataen. Joko silloin tai ei koskaan. Ja kaikki palasivat kotiin märkinä kuin koirat. Senvuoksi on jumalilla jalat villaan käärittyinä (s.o. he pitävät yllä lämmintä ja kuivaa ilmaa), että me emme enää ole hurskaita." (Trimalkion pidot, siv. 68—69.)

—16—

Kuva 24. Finereader 10 -ohjelman tekstintunnistusnäyte

Ohjelma oli selvästi paras vaihtoehto työn suoritukseen, mutta rajoitustensa (Taulukko 10, s. 33) takia ohjelmaa ei voitu käyttää.

Seuraavanlaisesti tulostui käännös (Kuva 25) PDF Converter 7 -ohjelmalla (k. 5.2.10, s. 40), Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat

virheet, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisia virheitä löytyi 6 kappaletta. Keltaisia virheitä löytyi 19 kappaletta.

vatus näin toisiinsa. Meidän katsomuksemme mukaan on yksilön kaikinpuolinen kehittäminen samalla parasta yhteisokasvatusta. Tämän mukaan on kasvatuksemme lahipaamaarana verevan suomalaisen ihmisen ja kunnan kansalaisen kehittäminen. Tätä **padmaarad** me pidämme koetinkivenamme kasvatuserojemme ja toimenpiteittemme oikeutusta ja arvoa **ntharattaessa**."

Kuten edellisestä ilmenee, tarkastelee maisteri Kaljunen artikkelissaan kasvatuskysymystä humanistisen etiikan pohjalta. Ihmisen uskonnolliset tarpeet, joita kirjoituksessa ei erikseen mainita, saavat talloin enintään tasa-arvoisina toisten tarpeiden rinnalla sijan kasvatuserojemmassa. Uskonnollisen elämän keskeinen merkitys persoonallisuuden muodostuksessa tulee siten kokonaan syrjäytetyksi

(vrt. 1 Tim. 4: 8: "jurnalaisuudesta on **hyotyä** kaikkeen, koska **sills** on **elarnan** lupaus, sekä nykyisen että tulevaisen"), puhumattakaan siitä menetyksestä, mikä kasvatuksen **paarnaria** on kohdannut niiden menetettyä tuonpuoleiset ikuisuustavoitteensa. Vrt. **mytis** sotilaspastori Raivion selostusta piispa Runestamin. kirjasta "Kristendom **oeh varlsforb.attring**" **lehternme Vasa**. numerossa.

Tässä yhteydessä ei voi olla viittaamatta siihen, **mink&** arvon jo pakanat ovat ymmärtäneet antaa uskonnolle kansan **ellimassa** ja **sen**. kasvatuksessa. Niinpa keisari Neron aikana elänyt **Petronius** Arbiter kirjoittaa kirjassansa Trimalkion pidot: "Mita tasta tuleekaan, jolleivat jumalat eivatka ihmiset saali **tats** kaupunkia! Niin totta kuin **tandon** saada iloa lapsistani, luulen, että kaikki **thma**. (nimittain tavaton kuivuus) tulee jumalilta. **Sill&** kukaan ei pida taivasta taivaana, ei kukaan noudata paastoa, ei kukaan valita **hittuistakaan** Juppiterista, vaan kaikki vain katsovat **alaspään** ja laskevat rahojaan.

Ennen aikaan menivät naidut naiset **pitkisa** vaipoissaan paljain jaloin vuorelle, hiukset hajallaan, mielet puhtaina, ja rukoilivat Juppiterilta vettä. Ja niinpa satoikin heti kuin saavista kaataen. Joko silloin tai ei koskaan. Ja kaikki palasivat kotiin **marking**, kuin koirat. Senvuoksi on jumalilla jalat villaan kaarittyina (s.o, he pitavat yllä lamminta ja kuivaa ilmaa), että **me** **ermine** **enad** ole hurskaita." (Trimalkion pidot, siv. 68-69.)

—16—

Kuva 25. PDF Converter 7 -ohjelman tekstintunnistusnäyte

Ohjelma oli virheineen selvityksen kelpuutettavin ohjelma, jonka rajoitukset eivät estäneet työn toteuttamista.

Seuraavanlaisesti tulostui käännös (Kuva 26) Adobe Acrobat pro 10 -ohjelmalla (k. 5.2.6, s. 35), Punaisella ovat virheet, joissa ohjelma on jättänyt sanasta tunnistamatta tai tunnistanut väärin enemmän kuin yhden merkin. Keltaisella ovat

virheet, joissa ohjelma on tunnistanut sanasta tai tunnistanut väärin vain yhden merkin. Punaisia virheitä löytyi 14 kappaletta. Keltaisia virheitä löytyi 25 kappaletta.

vatus nam toisiinsa. Meidän katsomuksemme mukaan on yksilön kaikinpuolinen kehittämisen samalla parasta yhteiskasvatusta. Tämän mukaan on kasvatuksemme lahja, am. aarana verevan suomalaisen ihmisen ja kunnan kansalaisen kehittämisen. Tata paamaaraii. me pidämme koetinkivenamme kasvatustilojemme ja toimenpiteittemme oikeutusta ja arvoa rrtiilirattaessa."

Kuten edellisestä ilmenee, tarkastelee maisteri Kaljunen artikkelissaan kasvatuskysymystä humanistis-etiikan pohjalta. Ihmisen uskonnolliset tarpeet, joita kirjoituksessa ei erikseen mainita, saavat tällöin e nintäin tasa-arvoisina toisten tarpeiden rinnalla -ajan kasvatustilassa. Uskonnollisen ehi. man keskeinen merkitys persoonallisuuden muodostuksessa tulee siten kokonaan syrjäytetyksi (vrt. 1 Tim. 4: 8: "jumalisuudesta on hyötyn kaikkiin, koska sillä on ehi. man lupaus, sekä nykyisen että tulevaisen"), puhumattakaan siitä menetyksestä, mikä kasvatuksen paamaaria on kohdannut niiden menetettyä tuonpuoleiset ikuisuustavoitteensa. Vrt. myöskin sotila. spastori Raivion selostusta piispa Runestamin kirjasta "Kristendom och varlsforbättring" lehtemme tässä numerossa.

Tässä yhteydessä ei voi olla viittaamatta siihen, minkä arvon jo pakanat ovat ymmärtäneet antaa uskonnolle kansan ehi. massassa ja sen kasvatuksessa. Niinpä keisari Neron aikana eliinyt Petronius Arbiter kirjoittaa kirjassansa Trimalkion pidot: "Mitä t. asta tuleekaan, jolleivät jumalat eiviitka ihmiset saali tata kaupunkia! Niin totta kuin tahdon saada iloa lapsistani, luulen, että kaikki tämä (nimittäin tavaton kuivuus) tulee jumalilta. Sillä kukaan ei pida taivasta taivana, ei kukaan noudata paastoa, ei kukaan valitii hi. mistakaan Juppiteristä, vaan kaikki vain katsovat alasp. ain ja laskvat rahojaan.

Ennen aikaan meniviit naidut naiset pitkissii vaipoissaan paljain jaloin vuorelle, hiukset hajallaan, mielet puhtaina, ja rukoilivat Juppiteriltä vetti. Ja niinpä satoikin heti kuin saavista kaataen. Joko silloin tai ei koskaan. Ja kaikki palasivat kotiin markkina kuin koirat. Senvuoksi on jumalilla jalat villaan kaiirittyyina (s.o. he pitaviit yllä lamminta ja kuivaa ilmaa), että me emme enää ole hurskaita." (Trimalkion pidot. siv. 68-69.)

- 1&-

Kuva 26. Adobe Acrobat pro 10 -ohjelman tekstintunnistusnäyte

Ohjelmaa ei virheiden määrän takia voitu valita selvityksen käyttökelpoisimmaksi.

Selvityksessä testatuista ohjelmista vertailtu Finereader 10 -ohjelma oli toimivin, mutta rajoitustensa (Taulukko 10, s. 33) takia ohjelman valitseminen työhön ei ollut mahdollista. Selvityksessä testatuista ohjelmista vertailtu PDF Con-

verter 7 -ohjelma oli puutteineen käyttökelpoinen, minkä takia ohjelma valittiin vuosikerran käsittelyyn.

7 Yhteenveto ja päätelmät

Työssä pääsin selvittämään, kuinka saada mahdollisimman tarkka digitaalinen versio käsin ladotusta kirjasta. Voi olla, että tulevaisuudessa tarvitaan tai halutaan tehdä toisesta kirjasta digitaalinen versio, josta on tarvetta hakea sanoja. Lopputuloksena sain PDF Converter 7 -ohjelmalla vuosikerran yhteen PDF-tiedostoon, johon pystyy tekemään sanahakuja. Ongelma lopputuloksessa oli harvaan kirjoitetut sanat, joita käytetyllä ohjelmalla oli vaikeuksia havaita. Vaikeuksia oli myös monen muun sanan kohdalla, mutta käsiteltäessä lähdetiedosto PDF-tiedostoksi, PDF Converter 7 oli käytettävissä olevista vaihtoehdoista käyttökelpoinen. Yritin selvittää, miten ohjelman saisi huomaamaan sanat, joita se ei tunnistanut tekstintunnistuksessa, mutta en saanut tähän kysymykseen vastausta. Word-muotoon muutettujen tiedostojen virheet on korjattava käsin, mutta parhaimmilla ohjelmilla käsiteltynä korjaukseen ei paljoa aikaa mene virheiden pienen määrän ansiosta. Haluttaessa löytää vuosikerrasta esimerkiksi luku 3, on hakukenttään kirjoitettava "N:o 3", jotta ohjelma löytää kyseisen luvun. Haluttaessa löytää luku päivämäärällä, on hakukenttään kirjoitettava esimerkiksi "28. 3. 46" siten, että pisteen jälkeen tulee väli.

Tietokantavaihtoehdon poissulkeminen oli ikävää, mutta yritys saada siitä toimiva olisi ollut liian suuritöinen ja ehkä loppujen lopuksi myös turha. Vaikka työn ohessa ilmeni omat ongelmansa, oli ohjelmien erilaisuus ja laatu mielenkiintoista selvitettävää.

Työasemalle ladattavien ohjelmien tulokset olivat pääasiassa parempia - eli niissä oli vähemmän virheitä - kuin selaimessa toimivat ohjelmat. Poikkeustapauksena oli selaimessa toimiva OCR Terminal -ohjelma, jonka lopputulos oli samaa tasoa työasemalle asennettavan Abbyy Finereader 10:n kanssa. Ilman rajoituksiaan (Taulukko 10, s. 33) eli rajattua käyttöaikaa ja rajattua tiedostonkäsittelymäärää Finereader 10 olisi ollut paras vaihtoehto projektin toteutukseen.

Testiaineistoista suurimmassa käytössä oli testiaineisto 2, mutta vaihtelevuuden ja kattavuuden saavuttamiseksi käytettiin myös neljää (4) eri testiaineistoa. Testiaineisto 2 valittiin pääkäyttöön sen monipuolisen sisällön takia. Monipuolisella sisällöllä tarkoitetaan eri tavalla kirjoitettuja sanoja. Testiaineistosta löytyi harvalla kirjasimella kirjoitettuja sanoja, eri tasossa olevia sanoja, lihavoituja sanoja ja myös numeroita.

Kuvat

- Kuva 1. Sotilasoppi –aikakausilehden ensimmäinen sivu, s. 10
- Kuva 2. Testausaineisto 2:n esittely (sivut 16 ja 17, jolla on luvun 2 alku), s. 11
- Kuva 3. Testauskohteen 3 esittely (sivut 2 ja 3 luvusta 1), s. 11
- Kuva 4. Testiaineisto 4:n esittely (sivut 50 ja 51 luvusta 4-5), s. 12
- Kuva 5. Testiaineisto 5:n esittely (sivut 32 ja 33, jolla on luvun 3 alku), s. 12
- Kuva 6. Free OCR –ohjelman tekstintunnistusnäyte, s. 19
- Kuva 7. Free OCR –ohjelman tekstintunnistusnäyte testiaineistosta 1, s. 20
- Kuva 8. Free online ocr –ohjelman tekstintunnistusnäyte, s. 21
- Kuva 9. i2OCR -ohjelman tekstintunnistusnäyte testiaineistosta 2, s. 23
- Kuva 10. New OCR –ohjelman tekstintunnistusnäyte, s. 24
- Kuva 11. New OCR –ohjelman tekstintunnistusnäyte testiaineistosta 2, s. 25
- Kuva 12. OCR Convert –ohjelman tekstintunnistusnäyte, s. 26
- Kuva 13. OCR Now! –ohjelman tekstintunnistusnäyte, s. 27
- Kuva 14. OCR Online –ohjelman tekstintunnistusnäyte, s. 28
- Kuva 15. OCR Terminal –ohjelman tekstintunnistusnäyte
testiaineistosta 4, s. 30
- Kuva 16. Merge PDF –selainpohjaisen ohjelman näkymä, s. 34
- Kuva 17. Adobe Acrobat pro 10 –ohjelman tekstintunnistusnäyte, s. 35
- Kuva 18. Autobahn DX:n .doc –version tekstintunnistusnäyte, s. 36
- Kuva 19. Autobahn DX:n PDF –version näkymä, s. 37
- Kuva 20. Autobahn DX –ohjelman näkymä, s. 37
- Kuva 21. Free OCR 3.0 –ohjelman tekstintunnistusnäyte, s. 38
- Kuva 22. Smart OCR –ohjelman tekstintunnistusnäyte, s. 39
- Kuva 23. PDF Converter 7 –ohjelman tekstintunnistus, s. 40
- Kuva 24. Finereader 10 -ohjelman tekstintunnistusnäyte, s. 42
- Kuva 25. PDF Converter 7 -ohjelman tekstintunnistusnäyte, s. 43
- Kuva 26. Adobe Acrobat pro 10 -ohjelman tekstintunnistusnäyte, s. 44

Taulukot

- Taulukko 1. Selainpohjaisten tekstintunnistusohjelmien kielet, s. 16
- Taulukko 2. Ohjelmissa käytetyt testiaineistot, s. 17
- Taulukko 3. Selainpohjaisten tekstintunnistusohjelmien käsiteltävät kuvaformaattit, s. 17
- Taulukko 4. Selainpohjaisten tekstintunnistusohjelmien tulostettavat tiedostomuodot, s. 18
- Taulukko 5. Selainpohjaisten tekstintunnistusohjelmien rajoitukset, s.18
- Taulukko 6. Työasemalle asennettavien tekstintunnistusohjelmien kielet, s. 31
- Taulukko 7. Ohjelmissa käytetyt testiaineistot, s. 32
- Taulukko 8. Työasemalle asennettavien tekstintunnistusohjelmien käsiteltävät kuvaformaattit, s. 32
- Taulukko 9. Työasemalle asennettavien tekstintunnistusohjelmien tulostettavat tiedostomuodot, s. 33
- Taulukko 10. Työasemalle asennettavien tekstintunnistusohjelmien rajoitukset, s. 33

Lähteet

1. Tekstien digitointi http://www.digiwiki.fi/fi/index.php?title=Tekstien_digitointi (luettu 2.5.2011)
2. Free OCR <http://www.free-ocr.com/> (luettu 3.5.2011)
3. Free Online OCR <http://www.free-online-ocr.com/> (luettu 3.5.2011)
4. i2OCR <http://www.sciweavers.org/free-online-ocr> (luettu 4.5.2011)
5. New OCR <http://www.newocr.com/> (luettu 4.5.2011)
6. OCR Convert <http://www.ocrconvert.com/> (luettu 5.5.2011)
7. OCR Now! <https://my.ocrnow.com/> (luettu 5.5.2011)
8. OCR Online <http://www.ocronline.com/> (luettu 6.5.2011)
- OCR terminal <https://www.ocrterminal.com/> (luettu 6.5.2011)
(myöhemmin:
9. <http://finereader.abbyyonline.com/en/Account/OcrTerminalWelcome>) (luettu 19.7.2011)
10. Abbyy <http://finereader.abbyy.com/> (luettu 11.7.2011)
11. Merge PDF <http://foxyutils.com/mergepdf/#> (luettu 14.7.2011)
12. Acrobat
https://www.adobe.com/cfusion/tdrc/index.cfm?product=acrobat_pro&loc=en
(luettu 14.7.2011)
13. Autobahn <http://www.aquaforest.com/en/autobahn.asp> (luettu 18.7.2011)
14. Free OCR http://www.freewarefiles.com/Free-OCR_program_34315.html
(luettu 19.7.2011)
15. Smart OCR <http://www.smartocr.com/> (luettu 20.7.2011)
16. Nuance PDF Converter <http://www.nuance.com/products/pdf-converter-professional7/index.htm> (luettu 10.8.2011)