

# POZNAŃ SUPERCOMPUTING AND NETWORKING CENTER



# POZNAŃ SUPERCOMPUTING AND NETWORKING CENTER



Country-scale infrastructure for creation of full  
text versions of historical documents from Polish  
Digital Libraries

Adam Dudczak, Miłosz Kmiecik, Marcin Werla  
{name.surname}@man.poznan.pl

Interedition Symposium  
Scholarly Digital Editions, Tools and Infrastructure  
19-20 March 2012

## Polish Digital Libraries Federation

- The Polish Digital Libraries Federation aggregates information from digital libraries distributed across entire Poland
- Operates since 2007
- Over 70 digital libraries
- Almost 920 000 metadata records



## Role of Polish DLF

- Common gateway to resources from Polish digital libraries
- Works as an intermediary between Polish digital libraries and Europeana
- Offers services for both regular users and digital librarians
  - Basic statistics
  - Map of Polish Digital Libraries
  - Duplicates detection
  - Digitisation plans
  - E-learning courses

## Resources in Polish DLF

- Almost 920 000 of digital objects, most of them are described as containing text
- Mostly newspapers, journals and books (72%)
- Oldest objects from 11<sup>th</sup> century,
- Most objects from 19<sup>th</sup> and beginning of 20<sup>th</sup> century
- DjVu is the most widely used format for representation of the text (79%)

## Description of the problem

- Having one common access point with metadata search is great but full text is better than most detailed metadata!
- Libraries contains a lot of materials which might be useful for researchers
- Lack of full text search results in poor visibility of those resources
- This issue is address by PSNC's work in the frame of SYNAT (<http://www.synat.pl>) project.

## How OCR is used?

- Survey was held in September-October of 2010
- We received responses from 26 major institutions
- Survey covered 70% of resources gathered by Polish DLF
- We were asking about all sort of things related to creation of full text search
  - Type and number of documents
  - Digitisation practices
  - Usage of OCR software

## Results of the survey

- Scanning resolution between 300-600 PPI, colour depth depends on type of resource
- FineReader and Document Express are most widely used OCR software packages
- Only 3 institutions use training capabilities of OCR software
  
- 40% of objects were a subject of OCR
- No one does the correction of OCR results



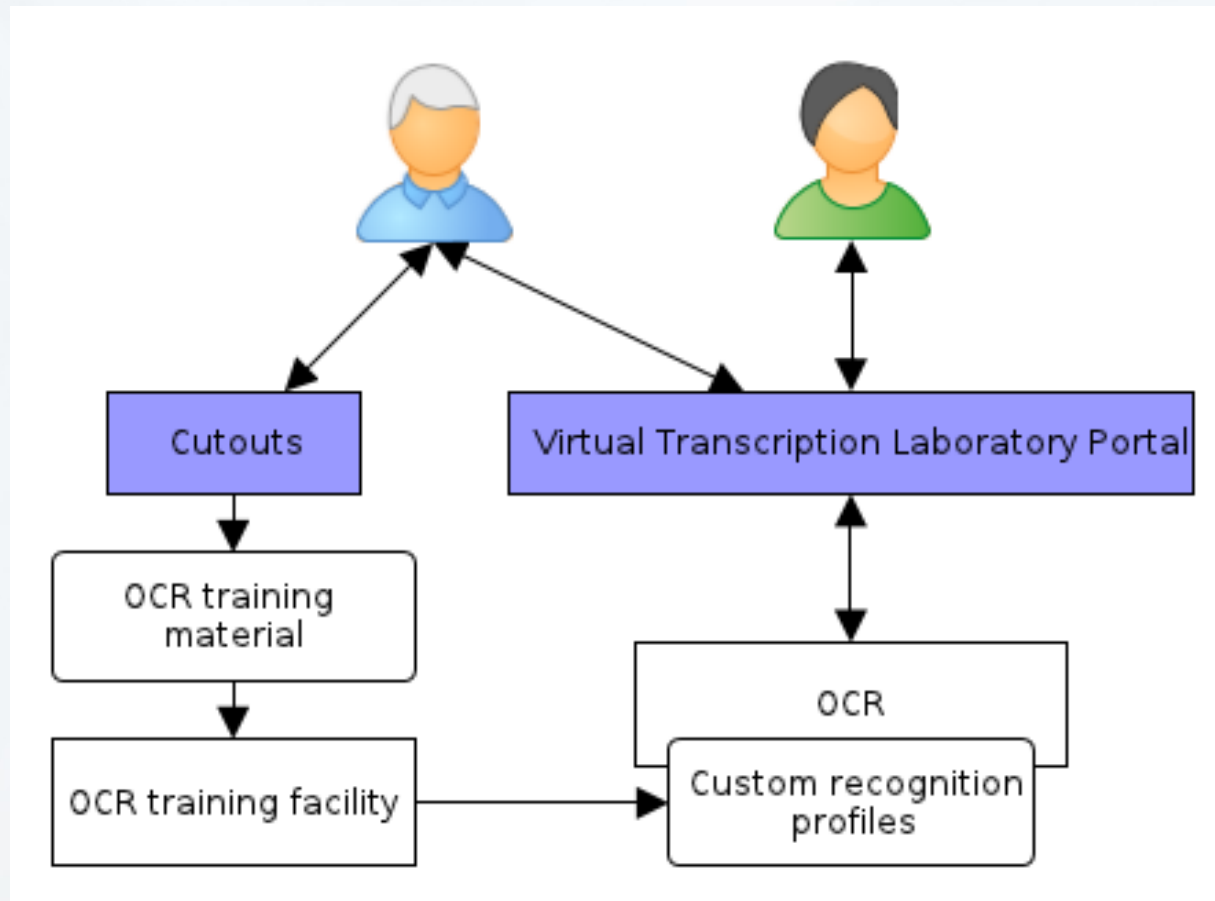
## Conclusions

- No correction of OCR results
  - Librarians are not interested in 100% correct text. OCR results are used as a search aid.
  - Lack of human resources to perform correction
- Limited usage of training capabilities
  - Training can improve OCR quality for historical documents
- Lack of tools which would integrate training and correction into digitisation workflow

## General assumption

- Integration with national aggregation infrastructure (Polish DLF)
- Well suited for digitisation workflow
  - Creation of custom recognition profiles
  - Crowdsourced correction of both new and existing resources
- Useful for researchers willing to work in distributed environment on historical projects

# Most important components

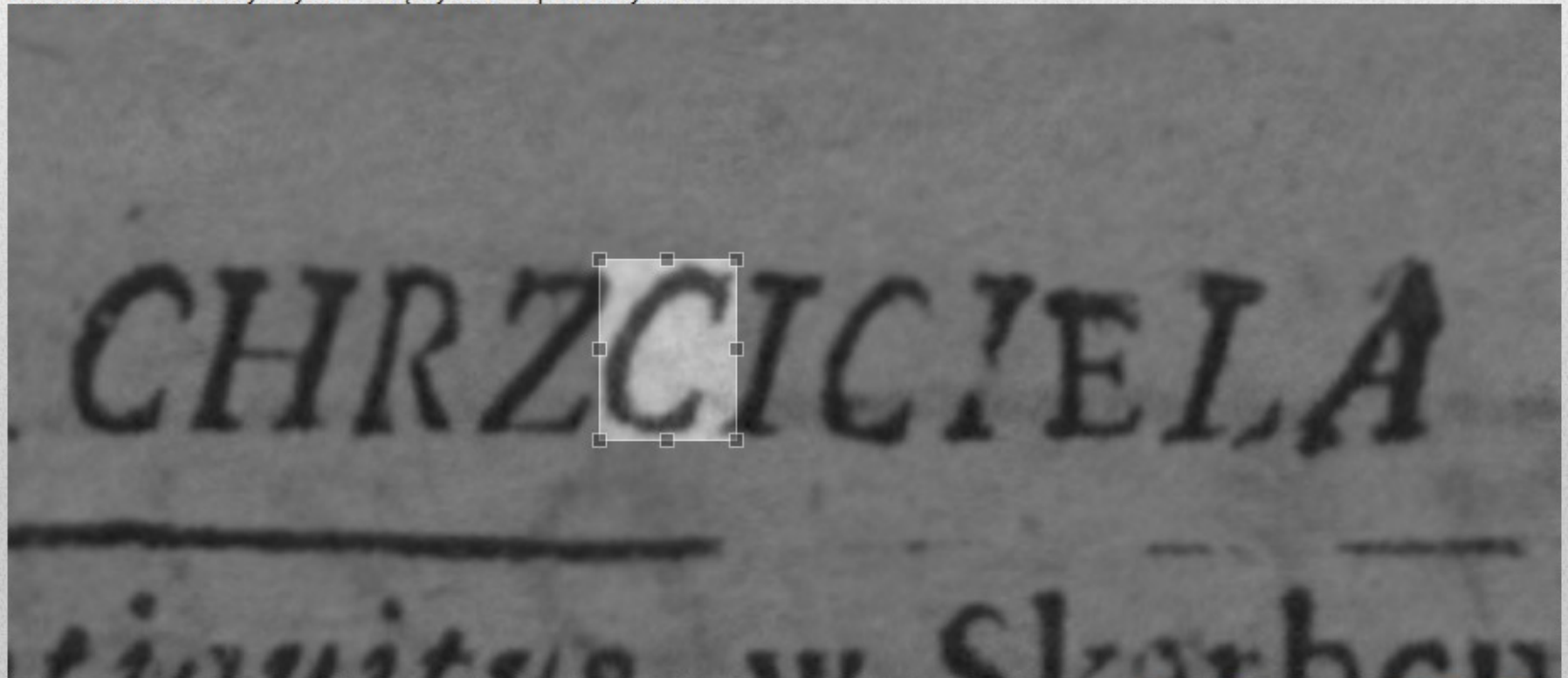


## OCR and supporting services

- OCR Service implemented on top of Tesseract 3.x
  - well-known, free, open source solution
- Support for recognition of multiple modern languages (including Polish)
- OCR supporting services
  - Preparation of training data (Cutouts, ...)
  - Training itself (OCR training facilities)

# Wycinanie znaków

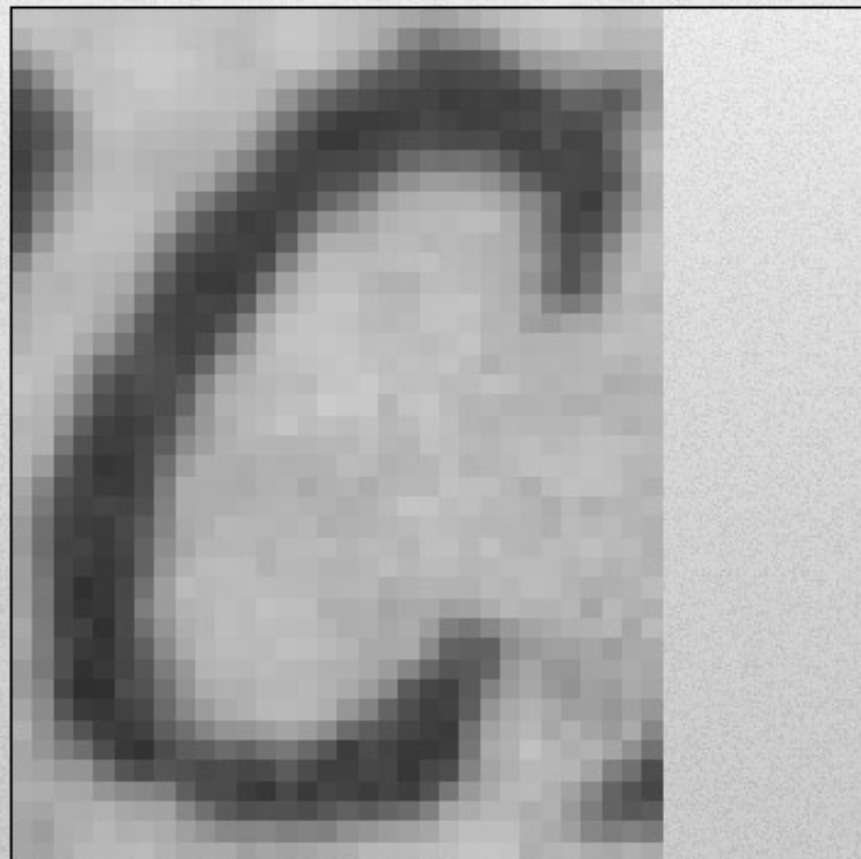
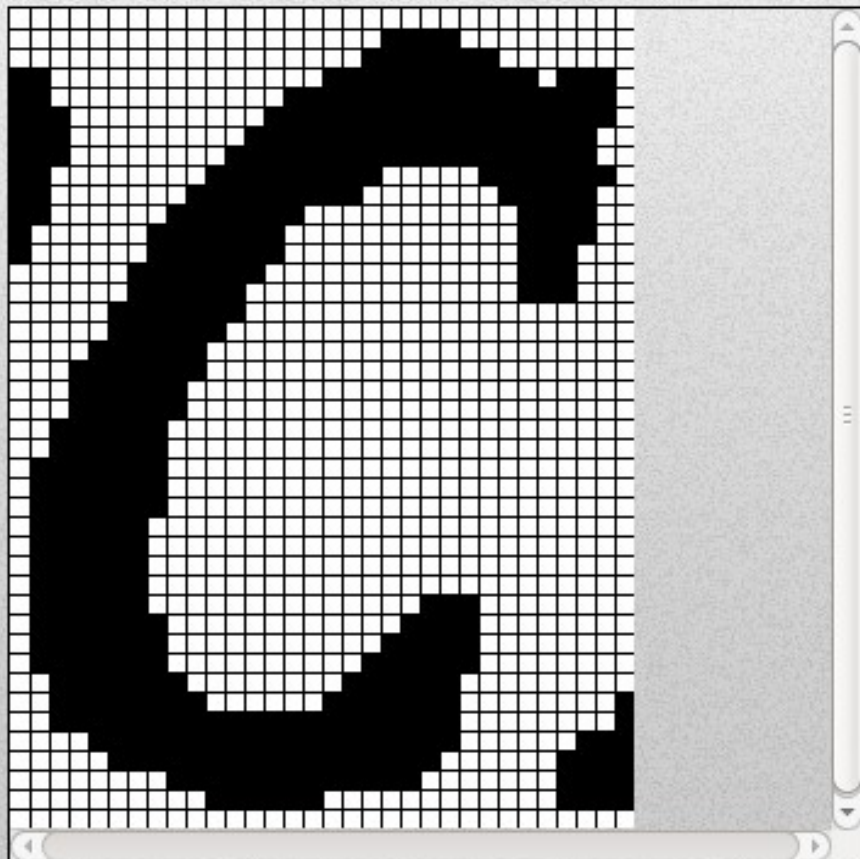
Pozostało: 2488 nieczytelnych: 1 błędnych: 7 opracowanych: 22



Historia opracowanych znaków		
znak	typ czcionki	opis
Z		<a href="#">wycofaj</a>
W		<a href="#">wycofaj</a>
a		<a href="#">wycofaj</a>
R		<a href="#">wycofaj</a>
R		<a href="#">wycofaj</a>
H	kursywa	<a href="#">wycofaj</a>

Screenshot from "Cutouts"

# Identyfikacja znaku



Zmień stopień binaryzacji - +



Korekty wykonane pędzlem [cofnij zmiany](#)

Formatowanie znaku:



Wpisz znak:

c



Screenshot from "Cutouts"

# Tesseract training

O A R C E P A N S K I E Ź 117

wfz yf ntencya, albowi m czytając *Adrychomiusza y Uni ersalem l alia Des-  
criptionem*, dochodzę, e w K ści le S. Jana Lat ran Źkim Rzymi , znay-  
dui się Arka Pańska, al ze złota al gołocona prz z Nie przyiaci .  
a si Rzymowi *eo od : itus y V sp ia* uszburzywłz Ognim y  
Mi cz m Jeruzalem, zabr wfz ni zliczone bogactwa, i *in triu pbo* do  
Rzymu w prow dzali, prowadzili y Arkę Pańską, lokując J w Kości l  
P koiu potym spalon m, t raz, wruder ch się pr z ntuiącym. Zaczym  
i k *Chr fti nitas* w Gł i Swi t swoje te podni a gł we; iak *erectum*  
*tropheu Crucis*, te Świę e *Lipsana* albo Reli wi Let raneńskiej dostały się  
Bazylic , oraz *Virga Aaron*, *Tabula* przykazania, *Virga Mo fis*, cz ty Kolu-  
mn z Jer zolimskiego Kościoła. Tam tedy y ta Święta Relikwia *asserua* ,  
al podobno ni ow ierwsza prz z Moyz sz uformowana, a Puszcz  
c . 15, lecz inna tamt y zrobionā, ktora odob o na go-  
rz N do tychczas utaiion ,  
Te y Arki Pańskiey, Lichtarza złot go, y Tytuśa Cezars na tr um-  
falnym Wozie siedzącego, na u al o i Tryumf lney są wyro-  
bion , ktory i Źnad Kościołem Nayswięt Panny Nowey w  
Rzymie. Zkąd t , w Tryumfie do Rz mu była Arka Pańska przy-  
pr adzona, owa powt rnie zrobiona, e y z tąd, ze według  
kiedy j lu rował Kościol Jerozolimski, iuż po w ściu ydow z Ba-  
bylonij y transport cyi Naczy ia z tamtąd, e widział Arkę, y a  
Autor takż H i świadczy,  
z Arkę w tryu fic prowadzono. , y inni Rzymu  
opif iący świadczą e , do tyc cza  
ieft Arka Pań a wyrżona, y na wozie Tryumfalnym siedzący,  
parā koni, y parā J dnorożcow iadący, prz z co, pi rwsze zdani dopie-  
ro namienione , i f z e lepiej *confirmatur*. Pod bno tedy, Źby tak po-  
ważny Autor w wyże allegowan ch e racy e nuiących  
ie refuuując ani , tak trzeba rozumieć, że Arka od Jeremiaśza  
na Gorze utaiiona, Niebu tylko wiadoma, pr y dokończeniu Świata  
swoie *effunde* skarby. A te o których piś g H storycy tylko yto sub-  
stytuowa c Arki przymierza, *ad instar* y na-  
k zt łt Rozgi Asonowej formo ane, ponieważ o nich w Piśmie S.  
Co zaś eolog Pański

O A R C E P A N S K I E Ź 117

wfzey ientencya, albowiem czytając *Adrychomiusza, y Unversalem Italia De-  
scriptionem*, dochodzę, że w Kościele S. Jana Lateraneńskim w Rzymie, znay-  
duie się Arka Pańska, ale ze złota cale ogołocona przez Nieprzyiaci .  
Dostała się Rzymowi *eo modo: Titus y Vespasianus* zburzywszy Ogniem y  
Mieczem Jeruzalem, zabrawszy niezliczone bogactwa, ie *in triumpho* do  
Rzymu w prowadzali, prowadzili y Arkę Pańską, lokując ją w Kościele  
Pokoiu potym spalonym, teraz, wruderach się prezentuiącym. Zaczym  
iak *Christianitas* w Głowie Świata swoje też podnieśła głowę; iak *erectum*  
*tropheum Crucis*, te Święte *Lipsana* albo Relikwie Leteraneńskiey dostały się  
Bazylicie, oraz *Virga Aaron*, *Tabula* przykazania, *Virga Moysis*, cztery Kolu-  
mny z jerozolimskiego Kościoła. Tam tedy y ta Święta Relikwia *asserua*,  
ale podobno nie owa pierwsza przez Moyześza uformowana, na puszczy  
*Exodi cap. 15*, lecz inna *ad formam* tamtey zrobionā, ktora podobno na go-  
rze *Nebo* do tychczas utaiiona,

Tedy Arki Pańskiey, Lichtarza złotego, y Tytuśa Cezars na tryum-  
falnym Wozie siedzącego, na *Arkusie* albo *Bramie* Tryumfalney są wyro-  
bione *vestigia*, ktory *Arkus* ieft nad Kościołem Nayswięt Panny Nowey w  
Rzymie. Zkąd *patet*, że w Tryumfie do Rzymu była Arka Pańska przy-  
prowadzona, owa powtornie zrobiona, *Patet* y z tąd, że według *Egesippa*  
kiedy *Pompejus* ułtrował Kościol Jerozolimski, iuż po wyjściu Żydow z Ba-  
bylonij y transportacyi Naczynia z tamtąd, że widział Arkę, y *Tabulas*  
*Testimonij*, item *Cherubinos*, Autor takż *Historia Scholastica Iudub* 3. świadczy,  
że Arkę w tryumfie Tytuśa prowadzono. *Adrichomius*, y inni Rzymu *spen-  
dorem* opifuiący *Authores* świadczą że *in Arcu triumphali Tytuśa*, do tychczas  
ieft Arka Pańska wyrżona, y *Tytus* na wozie Tryumfalnym siedzący,  
parā koni, y parā Jednorożcow iadący, przez co, pierwsze zdanie dopie-  
ro namienione, icż e lepiej *confirmatur*. Podobno tedy, Źby tak po-  
ważnych Autorow wyżej allegowanych e *probabiliter* racy onuiących  
nie refuuując ani *infringendo*, tak trzeba rozumieć, że Arka od Jeremiaśza  
na Gorze *Nebo* utaiiona, Niebu tylko wiadoma, prz dokończeniu Świata  
swoie *effunde* skarby. A te *Vasa* o których piś g Historycy tylko było sub-  
stytuowane *ad instar* Arki przymierza, *ad instar* *Tabularum Testimonij*, y na-  
k zt łt Rozgi Asonowej formowane, ponieważ *nulla* o nich w Piśmie S.  
*post Restaurationem Templi mentio* Co zaś Teolog Pański *Apoc: cap. 11. o*  
*Arce* temi wipomina słowy: *Et apertum est Templum DDI in Calo* *Et uisa est*  
*Arca*; nie ma się rozumieć, Źby Arka przymierza miała bydż do Nieba  
wzięta y przeniesiona, lecz przez nią rozumieć *Humanam Naturam Christi. Ki-  
sbardus*, y *Beda*, rozumieć drugdy Matkę Nayswiętszą, iako to *Bernardus* y  
*Suarez*,  
O Loto-

## Virtual Transcription Laboratory

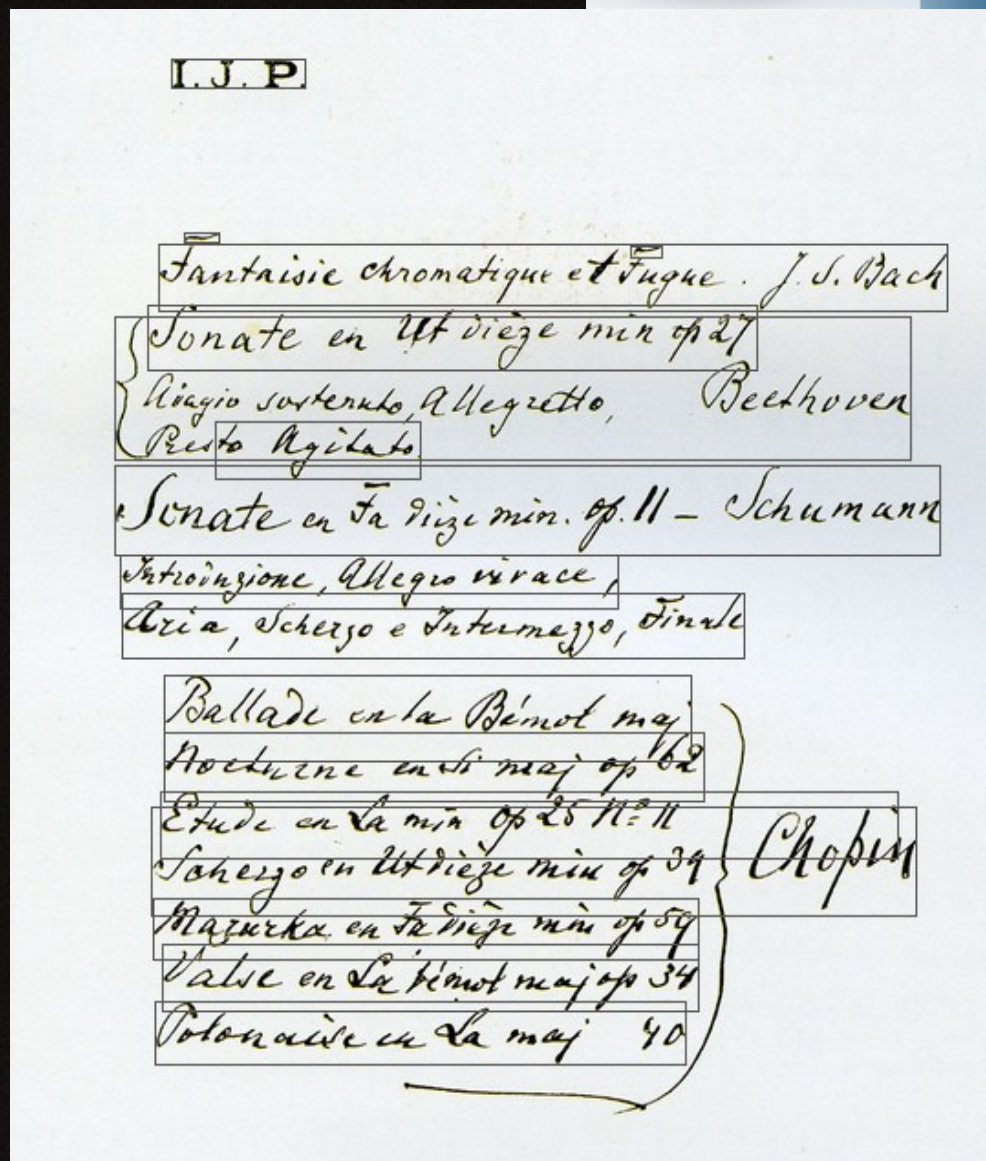
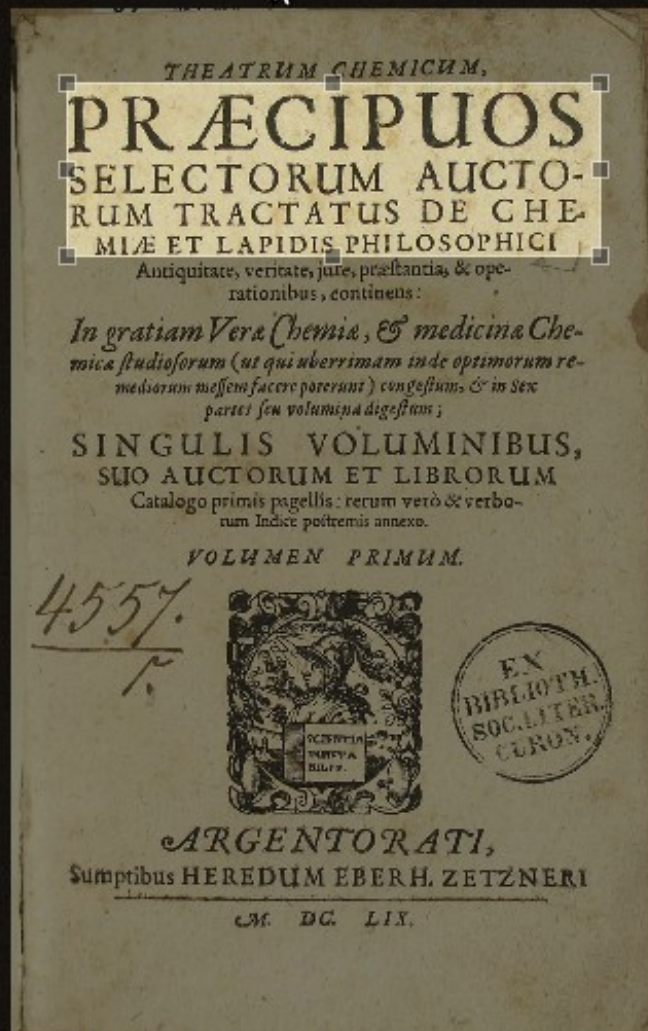
- Virtual Transcription Laboratory will allow to integrate text recognition and correction into digitisation workflows
- Users of VTL can upload scanned images and create textual version on this basis
- VTL gives access to OCR service
- It allows to correct existing text
- At any time user can export results of work in hOCR format



POWRÓT DO STRONY PROJEKTU

Zaznacz obszar, który chcesz poddać przetworzeniu OCR

-- default -- Start OCR



Create Transcript

Create Boundaries

tysiąc barek leżących w wodzie uniemożliwiało żeglugę. Odrzańscy wodniacy tylko dziesiątą część taboru rzeczno-objęli w posiadanie. Już w sierpniu 1945 roku jednak ruszył odrzańskim szlakiem pierwszy transport węgla.

Edytor

I

Komentarz

Line:

OCR

Wyczyść

Usuń



tysiąc barek leżących w wodzie uniemożliwiało żeglugę. Odrzańscy wodniacy

tego kanału aż do kopalń górnośląskich, aby uniknąć przeładowywania węgla na wagony kolejowe i z nich na barki. W roku 1960 „Żegluga na Odrze” przewiozła prawie 2 miliony ton towaru, za 5 lat przewiezie dwa razy tyle.

Odlóżmy notatki i nakreślmy na linii Odry te poprawki. Już nie ta rzeka, co piętnaście lat temu... Porównujmy dalej.

\*

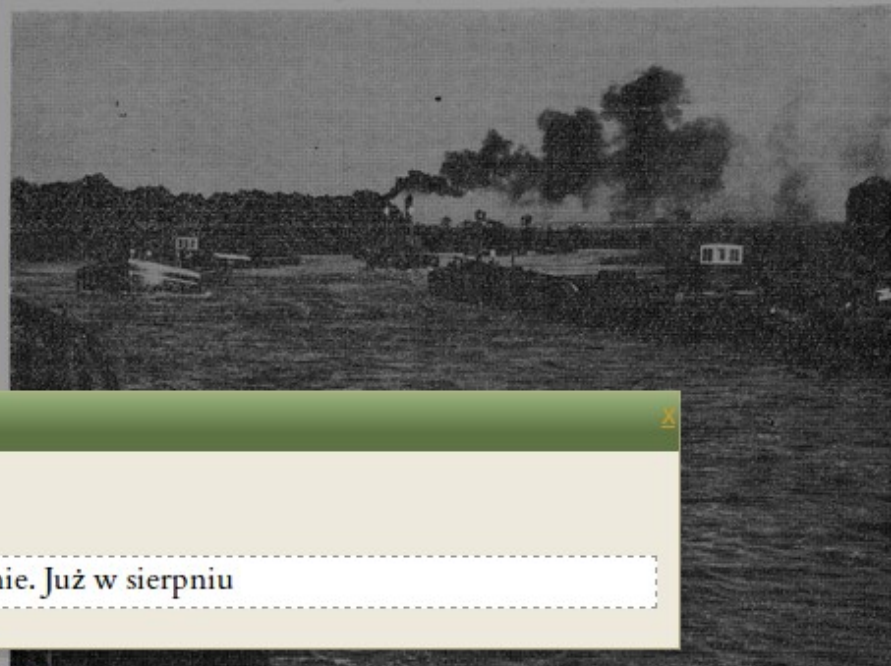
598 zakładów przemysłowych, skupiających ponad 63 tys. pracowników. Osiem wyższych uczelni. Ponad 15 tys. studentów. Dwa wydawnictwa... Wrocław.

Mówi się o nim, że jest stolicą Nadodrza i nie jest to tylko jakieś określenie honorowe. Wrocław jest potężnym ośrodkiem przemysłowym i kulturalnym, mającym niemały udział w kształtowaniu oblicza nie tylko Nadodrza, ale i całego kraju. Ludność Wrocławia stanowi 1,4% ludności Polski. Jego udział wynosi 2,6%. Zniszczenia wojenne są coraz mniej widoczne. Miasta przybywa około 6.000 izb mieszkalnych rocznie. Rosną nowe

# MIRTUALNE LABORATORIUM TRANSKRYPCJI



tysiąc barek leżących w wodzie uniemożliwiało żeglugę.  
Odrzańscy wodniacy  
tylko dziesiątą część taboru rzeczno-  
go objęli w posiadanie. Już w sierpniu  
1945 roku jednak ruszył odrzańskim szlakiem pierwszy transport węgla.  
Oczyszczono baseny portowe, odbudowano nabrzeża,  
zbudowano nowy stopień  
wodny, buduje się nowe zbiorniki retencyjne.



Edytor

I

Komentarz

Line: OCR

Wyczyść

Usuń



tylko dziesiątą część taboru rzeczno-  
go objęli w posiadanie. Już w sierpniu

Klawiatura

Poprzednie

Następne

č	Č	č	Ď	ď	Ď	ď	Ě	ě	Ě
ě	Ě	ě	Ě	ě	Ě	ě	Ĝ	ĝ	Ĝ
ğ	Ğ	ğ	Ğ	ğ	Ĥ	ĥ	Ħ	ħ	Ĩ
ĩ	Ĩ	ĩ	Ĩ	ĩ	Į	į	İ	ı	IJ
ij	Ĵ	ĵ	Ķ	ķ	κ	Ĺ	ĺ	Ļ	ļ
Ł	ł	Ł	ł	Ł	ł	Ń	ń	Ń	ņ

tysiąc barek leżących w wodzie uniemożliwiało żeglugę. Odrzańscy wodniacy  
tylko dziesiątą część taboru rzeczno-  
go objęli w posiadanie. Już w sierpniu  
1945 roku jednak ruszył odrzańskim szlakiem pierwszy transport węgla.  
Oczyszczono baseny portowe, odbudowano nabrzeża, zbudowano nowy stopień  
wodny, buduje się nowe zbiorniki retencyjne.

Przedsiębiorstwo „Żegluga na Odrze” otrzymało kredyty na budowę ponad 200 barek motorowych. Kanał Gliwicki łączy się obecnie magistralą wodną z Zakładami Azotowymi w Kędzierzynie. Trwają studia nad przedłużeniem tego kanału aż do kopalń górnośląskich, aby uniknąć przeladowywania węgla na wagony kolejowe i z nich na barki. W roku 1960 „Żegluga na Odrze” przewiozła prawie 2 miliony ton towaru, za 5 lat przewiezie dwa razy tyle.

Odlóżmy notatki i nakreślmy na linii Odry te poprawki. Już nie ta rzeka, co piętnaście lat temu... Porównujmy dalej.

598 zakładów przemysłowych, skupiających ponad 63 tys. pracowników. Osem wyższych uczelni. Ponad 15 tys. studentów. Dwa wydawnictwa... Wrocław.

Mówi się o nim, że jest stolicą Nadodrza i nie jest to tylko jakieś określenie honorowe. mającym n... i całego kra... w produkcji przemysłowej wynosi 20%. Zniszczona wojennie są coraz mniej widoczne. Miastu przybywa około 6.000 izb mieszkalnych rocznie. Rosną nowe

Screenshots from transcription editor

## Ongoing work

- Experimenting with Tesseract training on top of data released in IMPACT project
  - <http://dl.psnc.pl/activities/projekty/impact/results/>
- Release of custom recognition profile dedicated to Polish documents written in Gothic script
- Public release of Virtual Transcription Laboratory

## Future plans

- Direct import of content from digital library
- Inclusion of various crowdsourcing techniques e.g. games

# POZNAŃ SUPERCOMPUTING AND NETWORKING CENTER



QUESTIONS?

**Poznań Supercomputing and Networking Center**  
affiliated to the Institute of Bioorganic Chemistry of the Polish Academy of Sciences,  
ul. Noskowskiego 12/14, 61-704 Poznań, POLAND,  
Office: phone center: (+48 61) 858-20-00,  
fax: (+48 61) 852-59-54,

e-mail: [office@man.poznan.pl](mailto:office@man.poznan.pl), <http://www.man.poznan.pl>