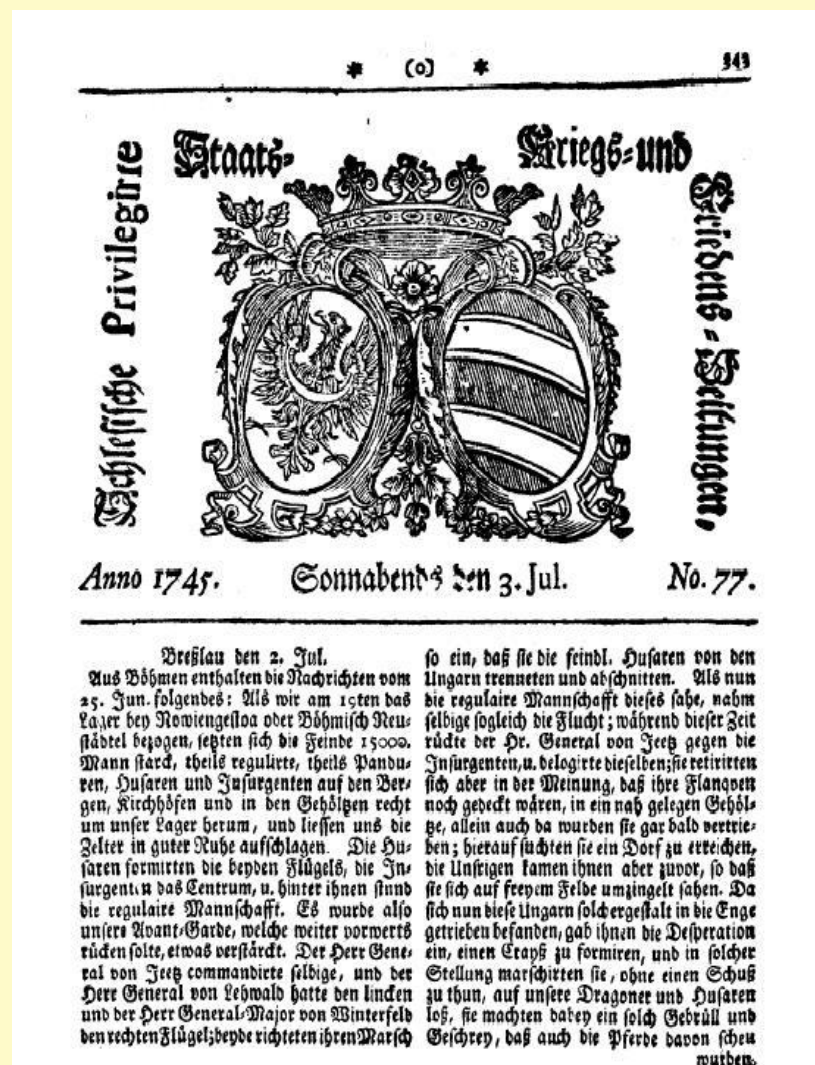
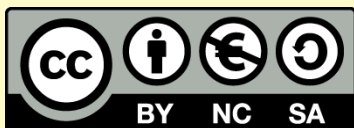


Korekta OCR – problemy i rozwiązania

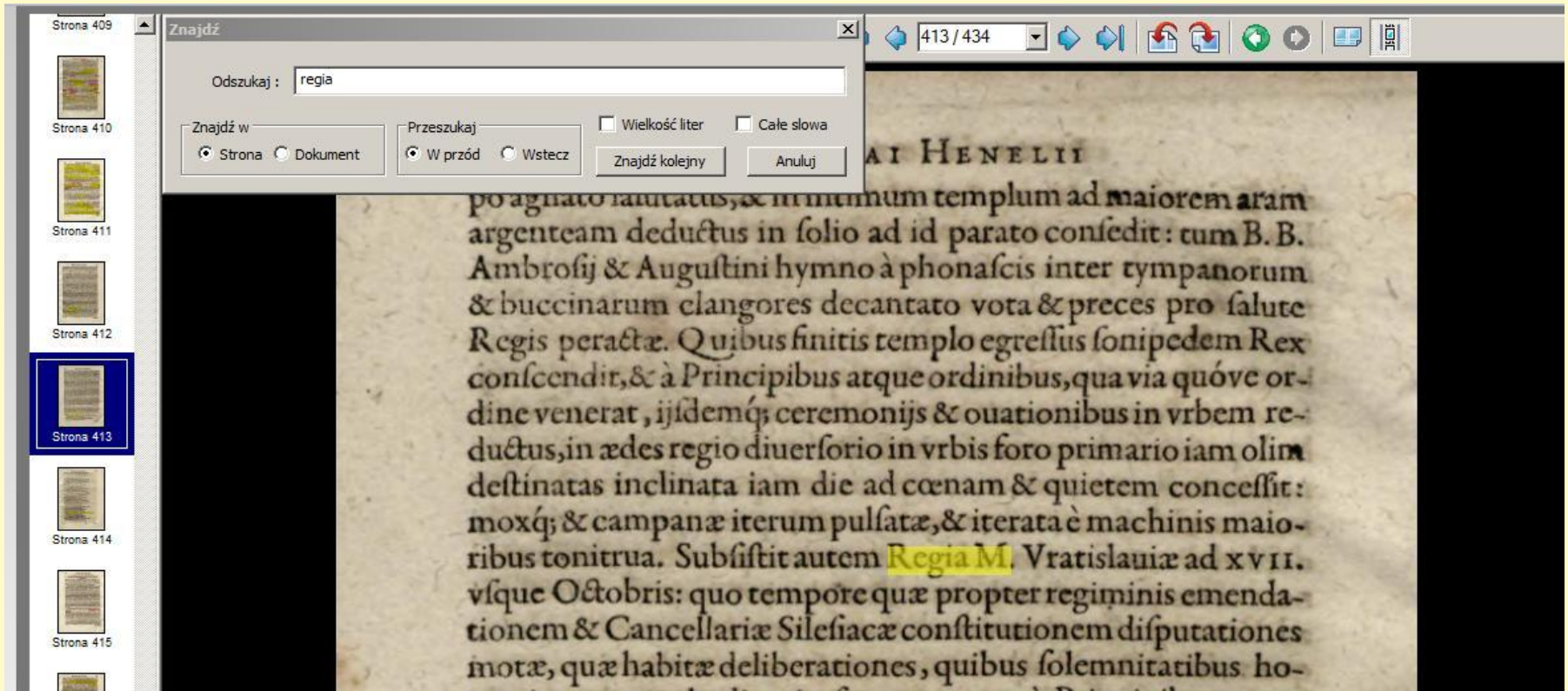
Edyta Kotyńska

eTEKA.com.pl



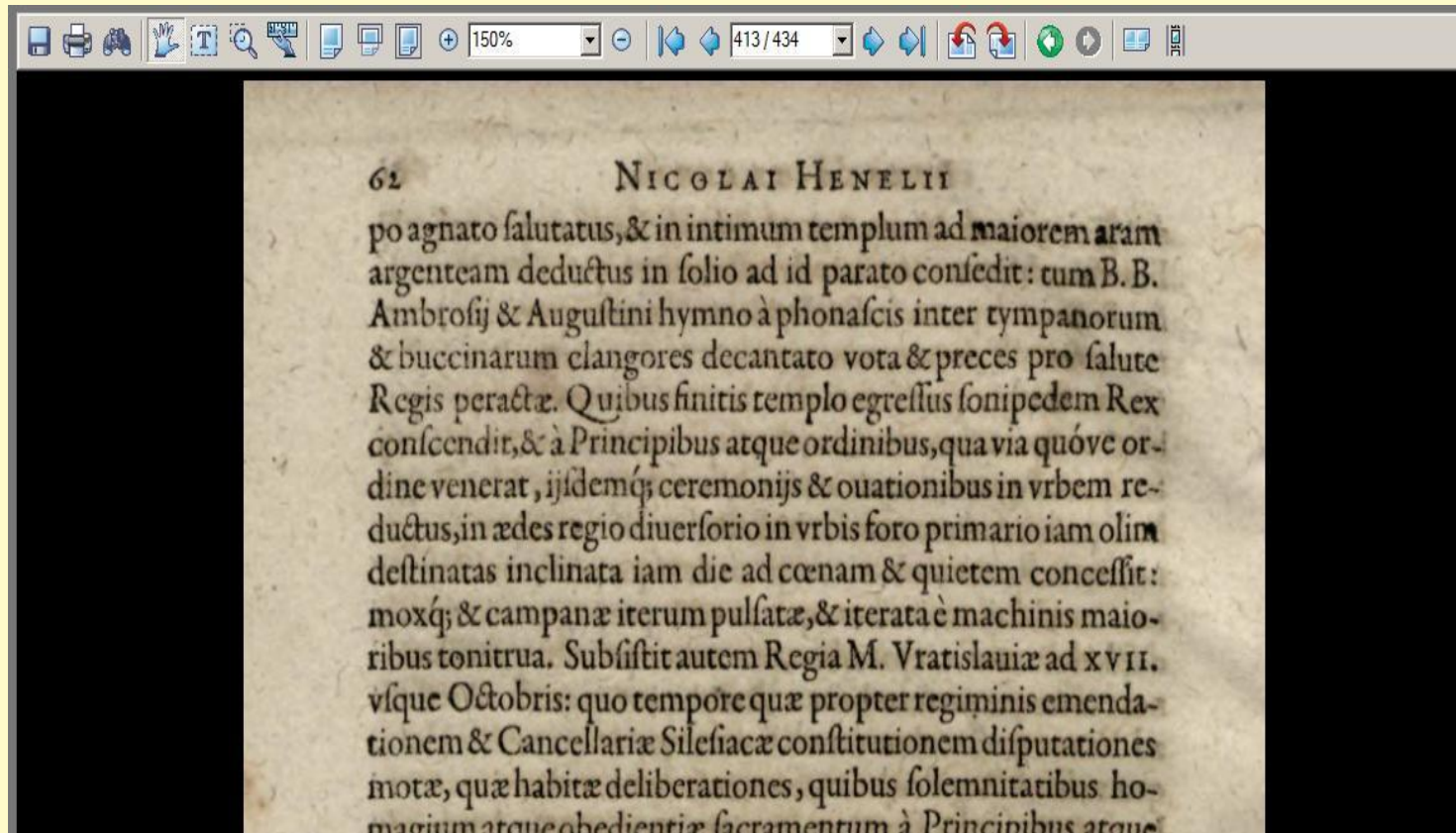
„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r.

Tekst, który nie jest cyfrowy - nie istnieje w sieci



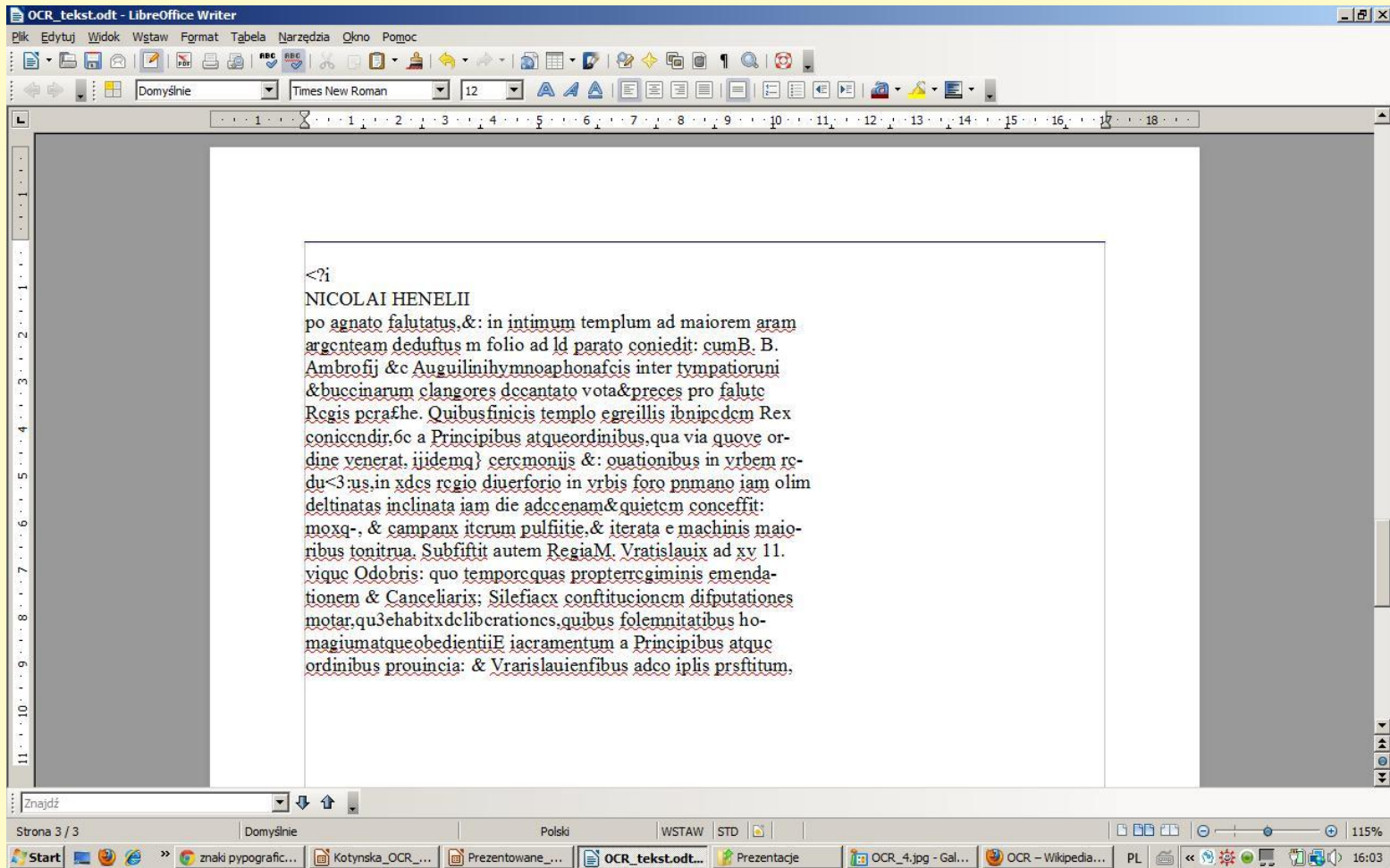
OCR umożliwia korzystanie z materiałów historycznych w formie przeznaczonej do odczytu komputerowego.

Materiały historyczne: obraz



„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

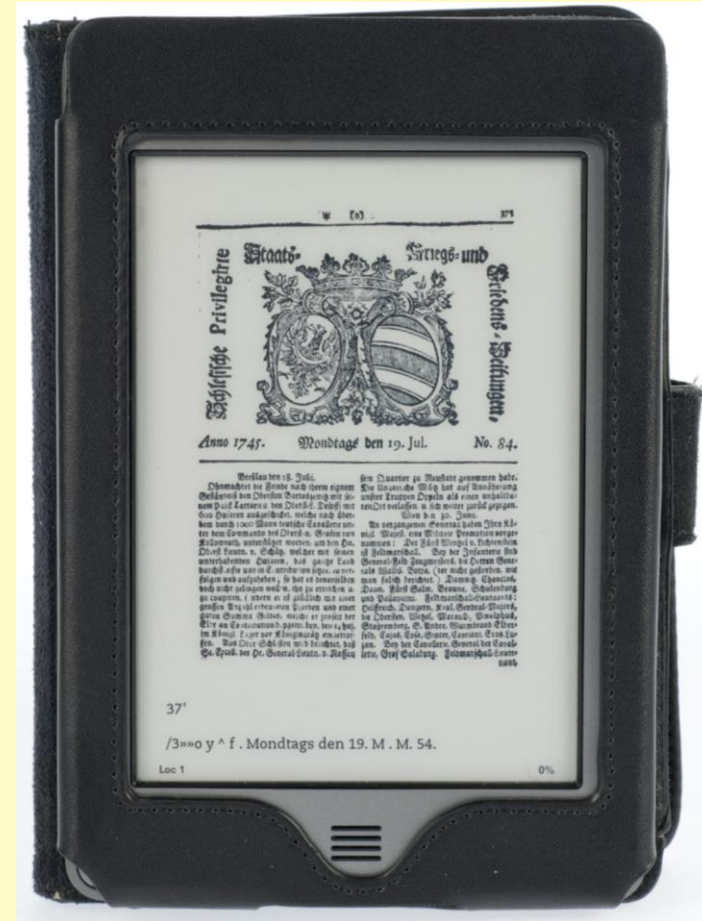
Materiały historyczne: tekst



„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

W efekcie OCR-u możemy uzyskać:

- pełnotekstowe indeksowanie i przeszukiwanie,
- możliwość kopiowania tekstu,
- powszechne udostępnienie materiałów „trudnych” do czytania,
- możliwość korzystania z translatorów,
- możliwość modyfikacji, uzupełnień,
- hiperlinkowanie w tekście,
- udostępnienie materiałów użytkownikom z dysfunkcją wzroku,
- umożliwienie wyeksportowania treści do formatów mobilnych.



Źródło: T. Kalota

OCR - optyczne rozpoznawanie znaków (ang. Optical Character Recognition) – zestaw technik lub oprogramowanie służące do rozpoznawania znaków oraz całych tekstów z obrazów cyfrowych i przetworzenie ich na teksty cyfrowe.

Zadaniem OCR-u jest rozpoznanie tekstu w zeskanowanym dokumencie:

- drukowanym,
- rękopiśmiennym,
- tabelarycznym.

ICR – zaawansowane rozpoznawanie kroju pisma, stopnia pisma, interlinii.

Terminy powiązane:

digitalizacja

zachowanie dziedzictwa kulturowego

dokumenty cyfrowe reprodukcja **obrazy cyfrowe**

teksty cyfrowe

badania naukowe **transkrypcja**

interpretacja ***biblioteki cyfrowe***

korekta konwersja słowniki

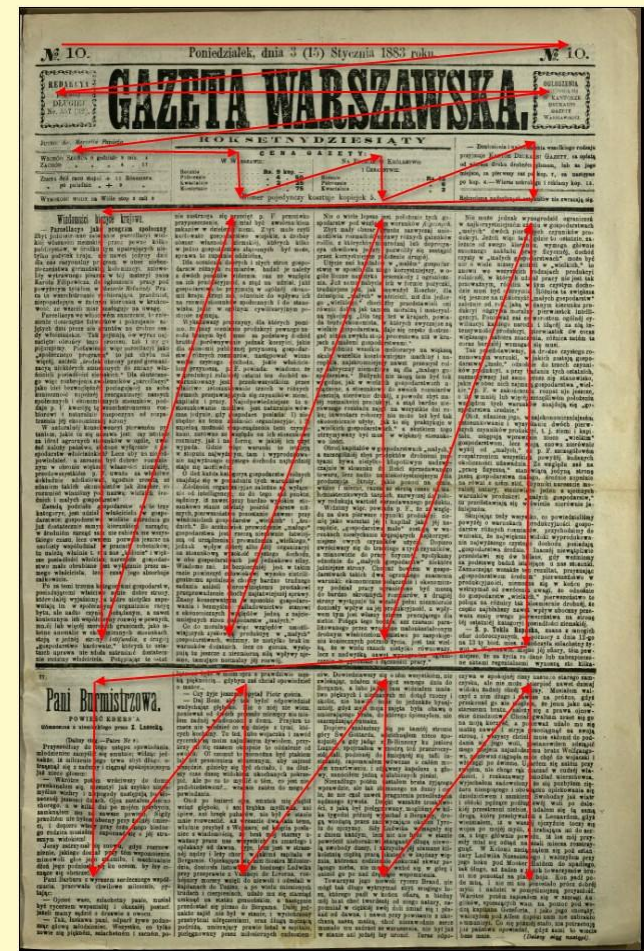
technologia

...



Dobry efekt rozpoznawania znaków zależy od:

- jakości tekstu i ogólnego stanu materiałów (zagięcia, pofałdowania),
- języka tekstu, alfabetu, kroju pisma,
- układu tekstu (kolumnny),
- formatowania tekstu (marginesy, interlinie),
- rozdzielczości skanowania i zastosowania głębi kolorów.

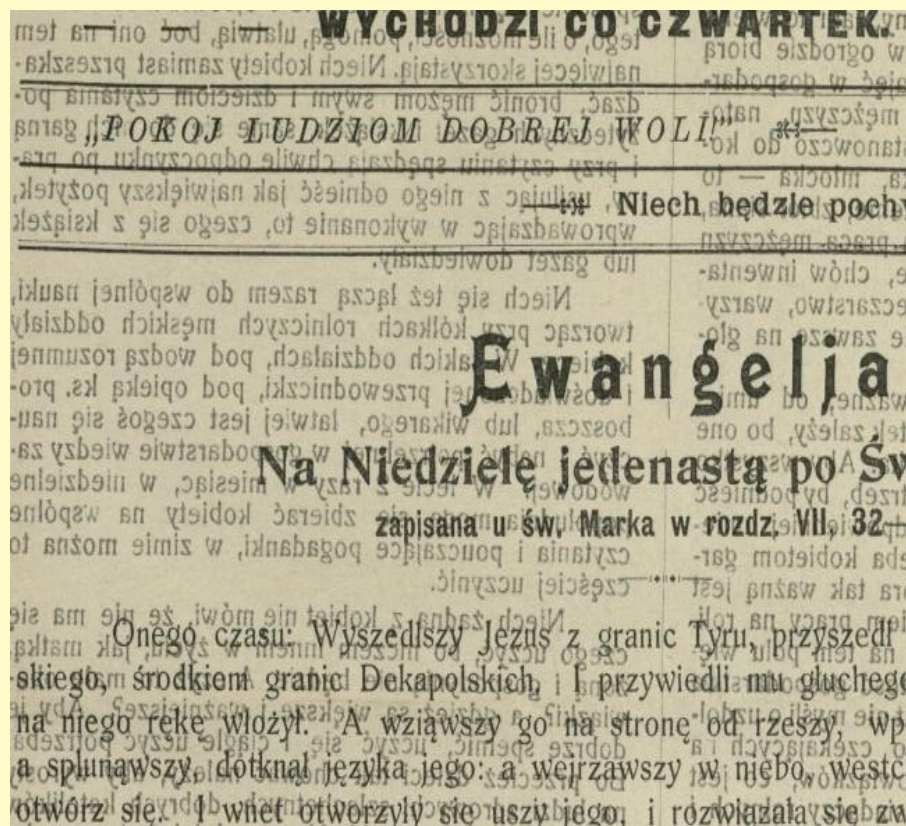


Źródło: G. Bednarek

Problemy techniczne dot. obrazu:

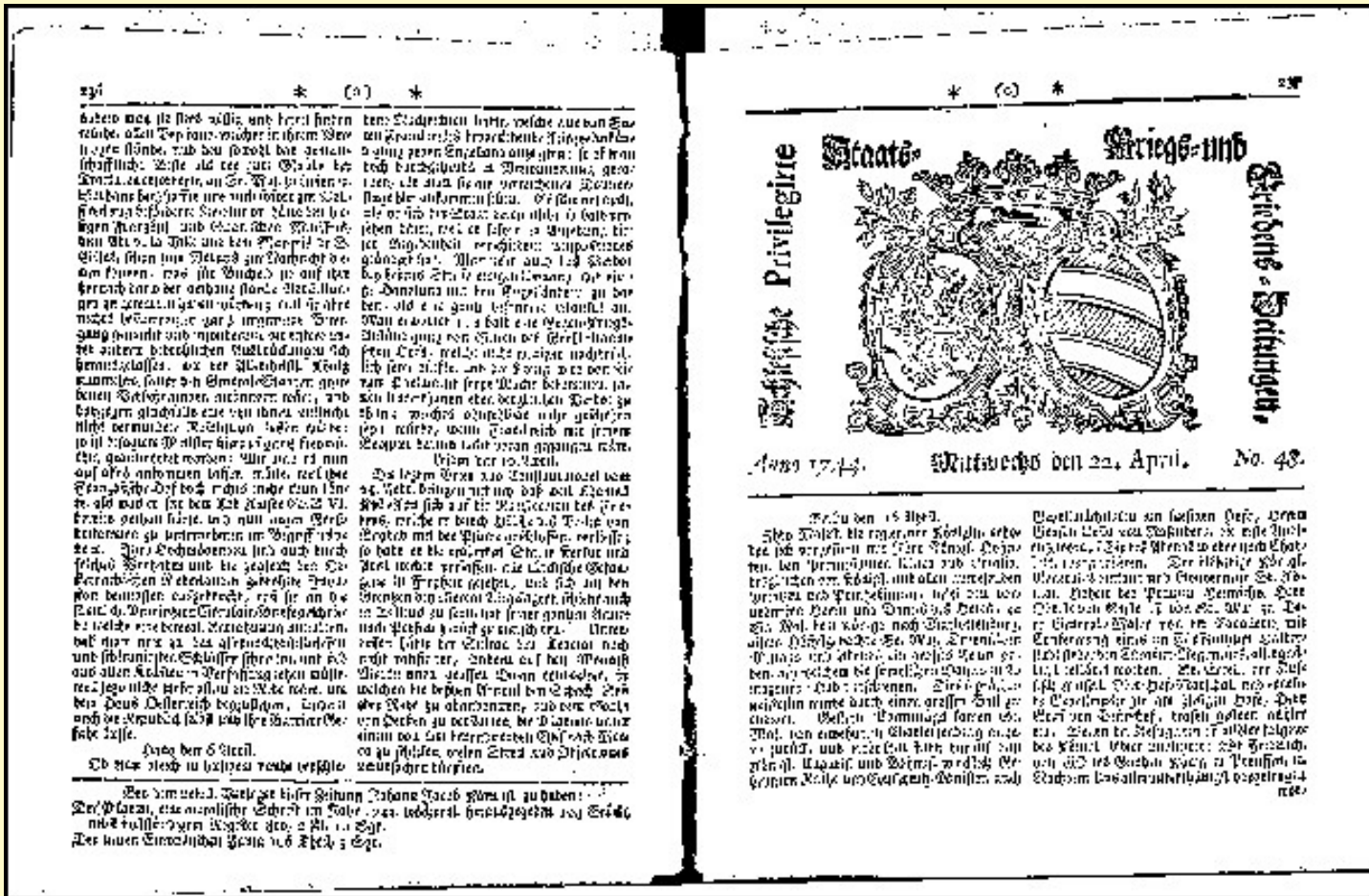
- wybór tekst (pomijanie grafiki),
- przebijający tekst,
- plamy, zalania, wandalizm,
- deformacja obiektu,
- kolorystyka i nasycenie druku,
- kolorystyka i nasycenie tła.

Konieczna obróbka techniczna.



Źródło: G. Bednarek

Plik przed obróbką



Źródło: T. Kalota



„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
 Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

Plik po obróbce

dabero man sie stets willig und bereit finden würde, allen Beystand, welcher in ihrem Vermögen stünde, und den sowohl das gemeinschaftliche Beste als der gute Glaube der Tractaten erforderte, an Sr. Maj. zu leisten. c. Sothane herzhafte und auch sofort zur Vollstreckung beförderte Resolution hätte den hiesigen Französis. und Spanischen Ministriß, dem Abt de la Ville und dem Marquis de S. Gillis, schon zum Voraus zur Nachricht dienen können, was für Bescheid sie auf ihre hernach darwider gethane starke Vorstellungen zu erwarten haben würden; weil sie aber nichts desto weniger ganz ungemene Bewegung gemacht und insonderheit der erstere unter andern bedrohlichen Ausdrückungen sich herausgelassen: wie der Allerchristl. König nunmehr, seiner den General-Staaten gegebenen Versicherungen entbunden wäre, und dazgegen gleichfalls eine von ihnen vielleicht nicht vermutete Resolution fassen würde; so ist besagtem Minister hierauf ganz freymüthig geantwortet worden: Wie man es nun auf alles antommen lassen müste, weil der Französische Hof doch nichts mehr thun könnte, als was er seit dem Tod Kaiser Carls VI. bereits gethan hätte, und nun gegen Großbritannien zu unternehmen im Begriff stünde. c. Ihre Hochmögenden sind auch durch solches Vorhaben und die zugleich den Oesterreichischen Niederlanden gedrohte Invasion vermassen aufgebracht, daß sie an die sämtliche Provinzen Circulaire-Briefe geschriebē, welche eine bewegl. Ermahnung enthalten, daß man nun zu den allergnädiglichsten und schleunigsten Schritten schreiten, und sich aus allen Kräften in Verfassung setzen müsse, weil jezo nicht mehr allein die Rede wäre, um dem Haus Oesterreich beizustehen, sondern auch die Republik selbst und ihre Barriere Gefahr lieffe.

Haag den 6 April.

Ob man gleich in hiesigem Lande verschiede-

Bei dem privit. Verleger dieser Zeitung Johann Jacob Korn ist zu haben: Der Pilgrim, eine moralische Schrift im Jahr 1743. wöchentl. herausgegeben 104 Stücke nebst vollständigem Register 8vo, 2 Fl. 10 Egr. Der neuen Europäischen Fama 106 Theil, 3 Egr.

dene Nachrichten hatte, welche eine von Seiten Frankreichs bevorstehende Kriegs-Ankündigung gegen Engelland anzeigten: so ist man doch durchgehends in Verwunderung gerathen, als man sie am verwichenen Donnerstage hier ankommen sehen. Es scheint auch, als ob sich der Staat deren nicht so bald versehen hätte, weil er sofort in Ansehung dieser Begebenheit verschiedene Dispositiones geändert hat. Man sieht auch das Verbot bey Lebens-Strafe einigen Umgang oder einige Handlung mit den Engelländern zu haben, als eine ganz besondere Clausul an. Man erwartet nun bald eine Gegen-Kriegs-Ankündigung, von Seiten des Großbritanischen Hofes, welche nicht weniger nachdrücklich seyn dürfte, und der König wird von seinem Parlament freye Macht bekommen, seinen Untertanen eben dergleichen Verbot zu thun; welches ohnfehlbar nicht geschehen seyn würde, wenn Frankreich mit seinem Beispiel darinn nicht voran gegangen wäre.

Leiden den 10. April.

Die letzten Brieue aus Constantinopel vom 25. Febr. bringen mit sich, daß weil Thamas Kuli-Kan sich auf die Ratification des Friedens, welche er durch Hülffe des Dacha von Bagdad mit der Pforte geschlossen, verliesse; so habe er die eroberten Städte Kerut und Aril wieder verlassen, alle türckische Gefangene in Freyheit gesetzt, und sich auf den Grenzen bey Mercat Alt gelagert, schiene auch in Willend zu seyn, mit seiner ganzen Armee nach Persien zurück zu marschiren. Unterdessen hätte der Sultan den Tractat noch nicht ratificiret, sondern auf den Monath Martii einen grossen Divan convociret, in welchen die beyden Articuli den Schach Esfi oder Kade zu abandoniren, und dem Sophi von Persien zu verstaten, die Pilgrims unter einem von ihm dependirenden Chef nach Mecca zu schicken, vielen Streit und Objectiones verursachen dürfften.



Anno 1744.

Mittwochs den 22. April.

No. 48.

Berlin den 16 April.
Ihre Majest. die regierende Königin, erheben sich vorgestern mit Ihrer Königl. Hoheit, den Prinzessinnen Ulrica und Amalia, desgleichen den Königl. und allen anwesenden Prinzen und Prinzessinnen, nebst den vornehmsten Herrn und Dames des Hofes, zu Sr. Maj. dem Könige nach Charlottenburg, alwo Höchstgedachte Sr. Maj. Denenfelden Mittags und Abends ein grosses Fein gaben, bey welchem die sämtlichen Dames in Russmazonen-Habit erschienen. Dieses prächtige Festin wurde durch einen grossen Ball geendiget. Gestern Vormittags kamen Sr. Maj. von erwehntem Charlottenburg abhero zurück, und ertheilten kurz darauf dem Königl. Ungarisch. und Böhmisch. würdlich Geheimen Rathe und Conferenz-Minister, auch

Genollmächtigten am hiesigen Hofe, Herrn Grafen Ustin von Rosenbergh, die erste Audienz, worauf Sie des Abends wieder nach Charlottenburg reiseten. Der bisherige Königl. General-Adjutant und Gouverneur Sr. Königl. Hoheit des Prinzen Heinrichs, Herr Obrist von Stille, ist von Sr. Maj. zu Hero General-Major von der Cavallerie, mit Conferirung eines im Fürstenthum Halberstadt stehenden Kürassier-Regiments, allergnädigt erklärt worden. Sr. Excell. der Russisch Kaiserl. Ober-Hof-Marschall und erwehnte Genollmächtige am hiesigen Hofe, Herr Graf von Destouches, trafen gestern alhier ein. Wegen der Defugirten ist alhier folgendes Königl. Edict publiciret: Wir Friedrich, von Gottes Gnaden König in Preussen, ic. Nachdem Uns allerunterthänigst vorgetragen

Źródło: T. Kalota

„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

Problemy techniczne dot. języka tekstu i alfabetu:

- łączenie słów,
- zamiana znaków,
- złe odczytywanie znaków,
- nieużywane współcześnie znaki i słowa.

ß ƒ Æ &

ss s AE et



Konieczna korekta

Problemy merytoryczne: obraz czy tekst?

- znaki specjalne,
- litery z akcentami,
- skróty (abrewiatury, ligatury),
- dawna pisownia nazw geograficznych i osobowych,
- dawna gramatyka i fleksja,
- zasób typograficzny drukarni,
- nieużywane współcześnie słowa.

<i>AE</i> → <i>Æ</i>	<i>ij</i> → <i>ij</i>
<i>ae</i> → <i>æ</i>	<i>st</i> → <i>ſt</i>
<i>OE</i> → <i>Œ</i>	<i>ft</i> → <i>ft</i>
<i>oe</i> → <i>œ</i>	<i>et</i> → <i>et</i>
<i>ff</i> → <i>ff</i>	<i>fs</i> → <i>ß</i>
<i>fi</i> → <i>fi</i>	<i>ffi</i> → <i>ffi</i>

Typowe ligatury alfabetu łacińskiego

SILESTOGRAPHIA.

fluuioli Brednicii haud procul Bithonia: & milliaria continet
ferme 20. Occiduum latus statuitur à fontibus Quissi fl. in
montibus Sudetis, ad confluentem vsque Boberi & Viadri,
qui est ad opidum Crofnam: & limes est Silesiæ Lusatiam ver-
sus & Marchiam Brandenburgensem, itidem milliarium 20.
& amplius. Austrinæ plagæ obiectum latus est à fontibus Vi-
stulæ vsque ad Quissi fontes complectens milliaria instar 40.
videlicet ab opido Tischena vsque ad fines Lusatiae. Pars se-
ptentrionalis à fontibus Brednicii ad confluentem vsque, qui
est Crofnæ: pariter milliaria circiter 40. continens. Regio
porro vniuersim tum montuosa est, tum campestris. Qua Se-

Plik Edytuj Widok Wstaw Format Tabela Narzędzia Okno Pomoc

Domyślnie Times New Roman 12

SiiEsroGRA* rftk,
 f
 ffutioli **Brednici** ihaudprocul **Bithonia**: &milliafiacontiftec
 fermezo. Occiduum lacus ftatuitur
 afontibus **Quissifl.** in **montibusSudetis**, ad tonfluentem vfque **Boberi**
 & **Viadri**,
 qui eftadopidum **Crosnam**: & limes eft **Silefiae Lufatiam** vet-
 ftis & **MarchiamBrandenburgensem**; itidem milhariurti 10.
 &ainplus. AuCtnnxglagxobic&um lacus eft afontibu^ **Vi-**
stulaevfque ad **Quissii** fontes complectens milliaria inftar 40»
 videlicet abomdo **Tischena** vfque ad fines **LufatiaJ**. Pars fe-
 pentnonais a fondbus **Brednicii** adconfluentem vfque, qui
 eft **Crosnae**: pariter milliaria ckciter 40. continens. Regio-
 pornavniuerfimumtuofaei^tumcampeftis. Qa.aSe-
 ptentrionem fpeftat & occidentem planior eft & abmonri-

Znajdź

Strona 1 / 3 Domyślnie Polski WSTAW STD 115%

po agnato salutatus, & in intimum templum ad maiorem aram
argenteam deductus in folio ad id parato confedit: tum B. B.
Ambrosij & Augustini hymno à phonaecis inter tympanorum
& buccinarum clangores decantato vota & preces pro salute
Regis peractæ. Quibus finitis templo egressus sonipedem Rex
conscendit, & à Principibus atque ordinibus, qua via quove or-
dine venerat, eisdemq; ceremonijs & ouationibus in urbem re-
ductus, in ædes regio diuersorio in urbis foro primario iam olim
destinatas inclinata iam die ad cœnam & quietem concessit:
moxq; & campanæ iterum pulsatae, & iterata è machinis maio-
ribus tonitrua. Subsistit autem Regia M. Vratislaviæ ad xvii.
vsque Octobris: quo tempore quæ propter regiminis emenda-
tionem & Cancellariæ Silesiacæ constitutionem disputationes
motæ, quæ habitæ deliberationes, quibus solemnitatibus ho-
magium atque obedientiæ sacramentum à Principibus atque

Heduuige czy **Hedwige**
Cvnradv czy **Cunradus**
Iablvncka czy **Jabluncka** czy **Jablunka**
Wwarta czy **Warta** czy **Wartha**

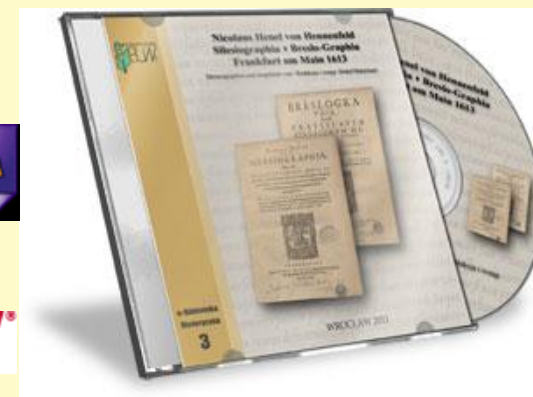
Korekta = Ekonomia

(sprzęt, oprogramowanie, licencje, praca, czas, efekty):

- korekta dla bibliotek cyfrowych,
- korekta dla projektów badawczych i wydawniczych,
- korekta społeczna,
- reCAPTCHA.



ABBYY

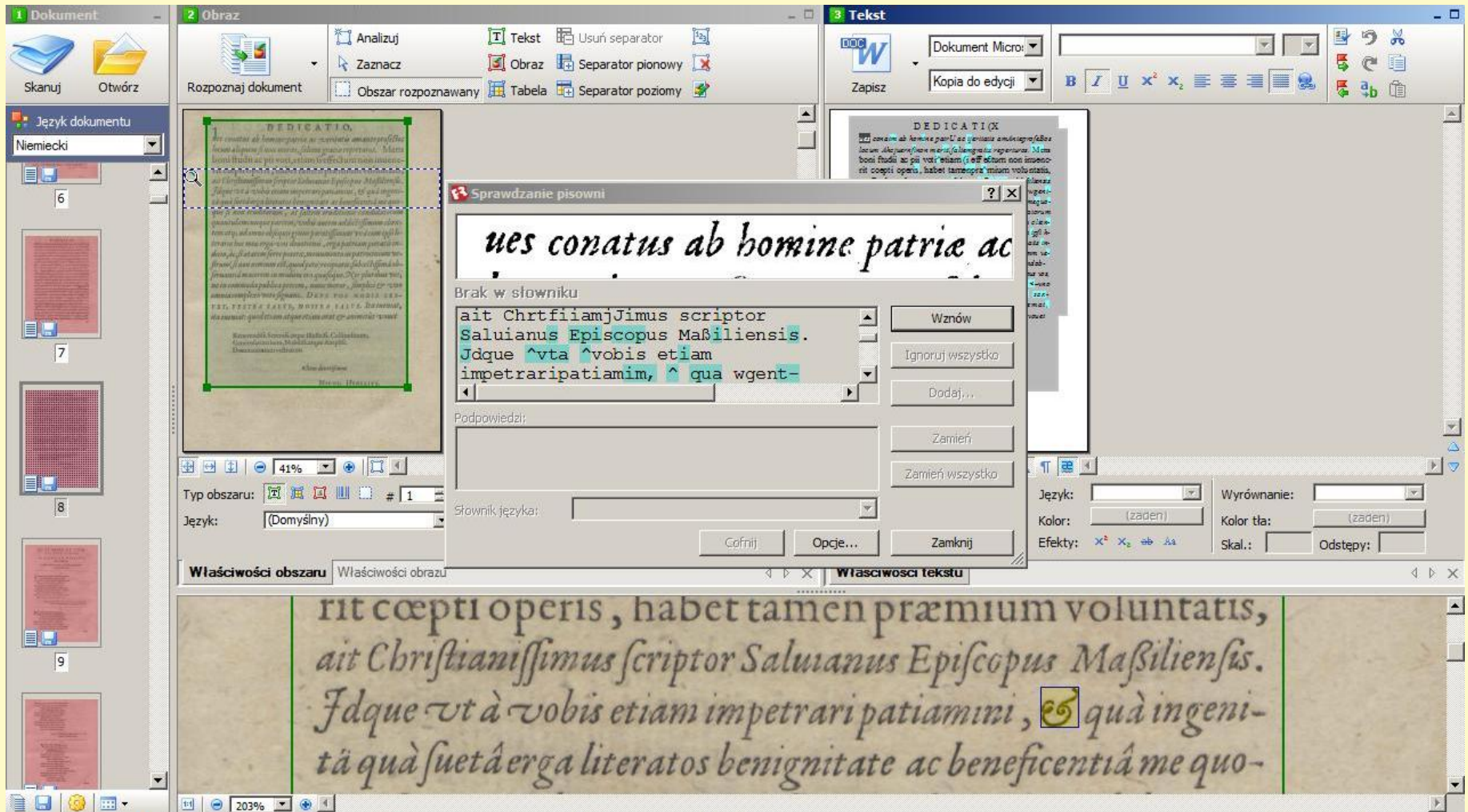


Project Gutenberg



„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

Korekta OCR Abbyy FineReader






„Biblioteka Cyfrowa dziś a wyzwania jutra” międzynarodowa konferencja naukowa
Kraków, 24-25.01.2013 r. Edyta Kotyńska: Korekta OCR – problemy i rozwiązania

Korekta OCR ReCAPTCHA

źródło: Technologie.org.pl

Tekst jest opublikowany na licencji Creative Commons – Użycie Niekomercyjne – Na tych samych warunkach 3.0.

Poleć:  Share  google +  e-mailem

Wyraż opinię: **podoba mi się 22** **nie podoba mi się 1**


SKOMENTUJ

Tytuł *:

Treść *:

Podpis *:

Adres e-mail:



Wpisz oba słowa:

Uwaga, komentarz pojawi się na liście dopiero po uzyskaniu akceptacji moderatora.

Autoryzacja komentarza pod artykułem <http://wiadomosci.ngo.pl/wiadomosci/833233.html>

WIRTUALNE LABORATORIUM TRANSKRYPCJI



STRONA GŁÓWNA

STWÓRZ PROJEKT

PROFIL

LOGOWANIE

KONTAKT



Co to jest?

Wirtualne Laboratorium Transkrypcji to portal, który wspiera digitalizację zasobów dziedzictwa kulturowego poprzez wspomaganie tworzenia



Jak to działa?

- Stwórz nowy projekt
- Wprowadź/zaimportuj skany
- Wykorzystaj naszą usługę OCR
- Transkrybuj/poprawiaj tekst dokumentu



Kto może korzystać z WLT?

- Naukowcy
- Instytucje kultury
- Hobbyści
- Użytkownicy bibliotek

Tweetnij 0

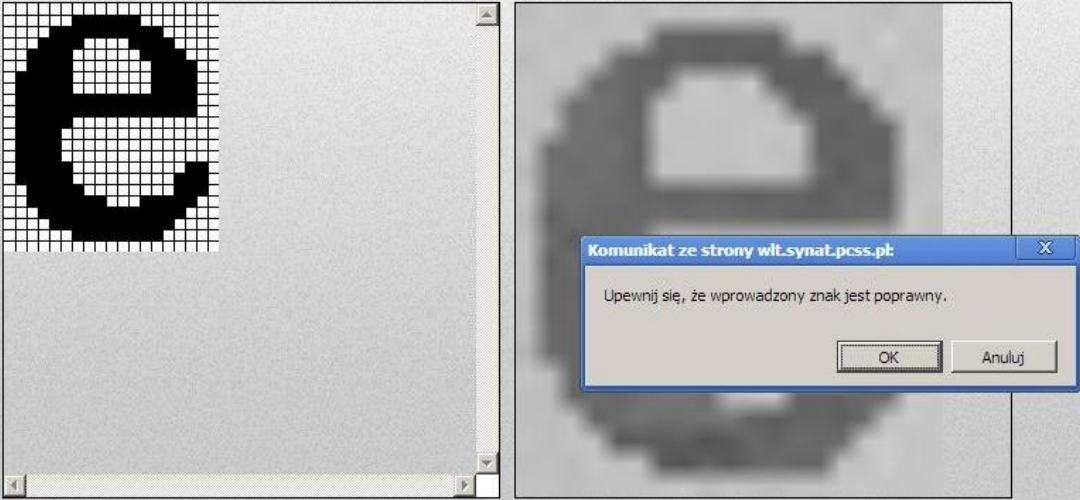
Like 3

Korekta znaków w WLT

Wycinanki przygotował Zespół bibliotek cyfrowych / PCSS

lista dokumentów | pomoc | kontakt | wyloguj

Identyfikacja znaku



Komunikat ze strony wlt.synal.pcss.pl

Upewnij się, że wprowadzony znak jest poprawny.

OK Anuluj

Wpisz znak: e

Dodatkowe cechy znaku: *I* **B** U

Zmień stopień binaryzacji - +

Korekty wykonane pędzlem cofnij zmiany

Projekt strony na podstawie szablonu autorstwa [NodeThirtyThree](#)

Podsumowanie:

Cel – dążenie do pełnej użyteczności historycznych materiałów, które są już dostępne w sieci.

OCR jest nadal w fazie bardzo intensywnego rozwoju z powodu dobrych narzędzi, ale ciągle jeszcze niedoskonałych efektów ich pracy. Korekta OCR-u wykonywana przez człowieka jest efektywna, ale bardzo kosztowna. Masowa digitalizacji, aby była opłacalna wymaga automatyzacji procesów, a tym samym ciągle doskonalszych narzędzi wspomagających i „uczących” OCR. Dodatkowo generalnych ustaleń wymagają problemy merytoryczne.

Wybrane źródła:

Repozytorium Instytucjonalne PCSS: <http://lib.psnc.pl/dlibra>

Biblioteka 2.0: <http://forum.biblioteka20.pl/index.php>

Digitalizacja.pl: <http://digitalizacja.pl/>

Format Djvu: <http://www.djvu.com.pl/>

Federacja Bibliotek Cyfrowych: <http://fbc.pionier.net.pl/owoc>

Wybrani autorzy publikacji dot. OCR-u:

Grzegorz Bednarek, Adam Dudczak, Tomasz Kalota, Tomasz Parkoła, Paweł Rękar, Marcin Szala



Anno 1745.

Sonntags den 3. Jul.

No. 77.

Breslau den 2. Jul.
Aus Böhmen enthalten die Nachrichten vom 25. Jun. folgendes: Als wir am 15ten das Lager bey Nowiengenloa oder Böhmisches Neustädtel bezogen, setzten sich die Feinde 15000. Mann stark, theils regulirte, theils Panduren, Husaren und Insurgenten auf den Bergen, Kirchhöfen und in den Gehöften recht um unser Lager herum, und lieffen uns die Zelter in guter Ruhe aufschlagen. Die Husaren formirten die beyden Flügel, die Insurgenten u das Centrum, u. hinter ihnen stand die regulaire Mannschafft. Es wurde also unsere Koant-Garde, welche weiter vorwärts rücken sollte, etwas verstärkt. Der Herr General von Zeeb commandirte selbige, und der Herr General von Lehwald hatte den linken und der Herr General-Major von Winterfeld den rechten Flügel; beyde richteten ihren Marsch

so ein, daß sie die feindl. Husaren von den Ungarn trenneten und abschnitten. Als nun die regulaire Mannschafft dieses sah, nahm selbige sogleich die Flucht; während dieser Zeit rückte der Hr. General von Zeeb gegen die Insurgenten, u. belogirte dieselben; sie retirirten sich aber in der Meinung, daß ihre Flancken noch gedeckt wären, in ein nah gelegen Gehölge, allein auch da wurden sie gar bald vertrieben; hierauf suchten sie ein Dorf zu erreichen, die Unstigen kamen ihnen aber zuvor, so daß sie sich auf freyem Felde umzingelt sahen. Da sich nun diese Ungarn solckergestalt in die Enge getrieben befanden, gab ihnen die Desperation ein, einen Crayß zu formiren, und in solcher Stellung marschirten sie, ohne einen Schuß zu thun, auf unsere Dragonen und Husaren los, sie machten dabey ein solch Gebrüll und Geschrey, daß auch die Pferde davon scheu wurden,

Dziękuję za uwagę. Edyta Kotyńska

eTEKA.com.pl

edyta.kotynska@eteka.com.pl

Digitalizacja.pl

