

The Fourth SEEDI Conference
Digitization of cultural and scientific heritage
Belgrade, Serbia, June 12-15, 2008

A software tool for searching in binary text images

Nikolay Kirov Kirov

Computer Science Department, NBU and
Institute of Mathematics and Informatics, BAS
Sofia, Bulgaria

Abstract: In this paper we present a software tool for searching word images in scanned text documents. We consider that the document pages are represented as files in tif, jpg, gif, png, bmp and other graphic file formats. Our experiments prove the efficiency of the proposed approach and show that such type of search could be successful. Examples of using various languages are presented. The software is user oriented and can be applied to any collection of scanned documents.

Introduction

Optical character recognition (OCR) is the usual way of conducting text retrieval from scanned document images.

It converts text images into a text file, recognizing every symbol and mapping it to a number, which is called code.

The most often used codes are ASCII (one byte code) or UTF-8 (two bytes code).

This technique is well developed and has high accuracy.

Searching words in a text file is relatively easy task.

But sometimes OCR is a very difficult process requiring dictionaries in the corresponding languages. Often human efforts are needed to correct OCR errors which is quite tedious work. There are some obstacles to successful OCR:

- The quality of page images.
- Language dependency (alphabet and coding, unknown language):
 - dictionaries;
 - old grammar, obsolete words and phrases and idioms;
 - old letters, outside of the coding tables;
 - multi-lingual documents;
- Errors in automatic OCR, human intervention needed.

de la Terre, &c.

fort peu de sens commun.

Quant à la Terre, si vous la **rencontrez** bonne, ce vous sera un grand avantage, & une grande épargne ; mais rarement en pourrez-vous trouver, où il n'y ait beaucoup à travailler, d'autant que telle paraîtra passablement bonne au dessus, qui étant ouverte de la profondeur d'un fer de Béche seulement, se trouvera Argileuse dessous ; ce fonds est pire aux Arbres que le Tuf, ou la Roche, à cause qu'il s'y **rencontre** de petites veines où les Racines peuvent s'étendre & profonder, afin de tirer la fraîcheur de plus bas, & prendre quelque nourriture ; mais l'Argileuse ou Terre franche ou rouge, fait comme un plancher qui par sa dureté & densité, ne peut être percé par aucunes Racines, & qui dans les grandes ardeurs de l'Eté, em-

il)uanrГ¤, la l'erre, ^! vous la **rencomrex** t,onne, ce vous lera un ^ranГ¤ avanra^e, sc u/B»e ^ranГ¤B« epar^ne ; maiz raremenr en pourrex-vouz trouvex, on il n've air tieaucouri a travailler, claur^nc ^ue teile riarrrra rrallГ¤llernenr cionne au clelluz, c^am ouverre cle la rirokoncleur Г¤'un iec Г¤e Lecrre ieulemenr, re rrouveral ^r?ileule cleclouz ; ce fonclz eli riire aux ^rГ¤rez c^ue le l"uf, on la l^t,cke, a caule c^u'il z'^ **reunorre** cle r/erirez veuez on lez ^ .acinez peuvenr z'ccenГ¤resc rirofoncler, arlB« cle rirer la ir, uclleurcle^luz baz, sc rirenГ¤re ^uelc^ue nourricure; malz l^r^ilcule ou ' I'errre francrre ou rou^e, fair comme un plancrer c^ui riar ia Г¤urere sc clenirre, ne peur ecre rierce riar aucunez li.acinez, sc c^ui cl^nz lez ^r, inclez arcleurz cle l'^, te, emB»

We suggest a different approach: instead of applying two steps – OCR and searching in text documents, we will directly search words in scanned text documents.

We can organize retrieval of words, similar to a given pattern word, searching in the binary text images.

The document pages can be represented as binary images in any graphic file format.

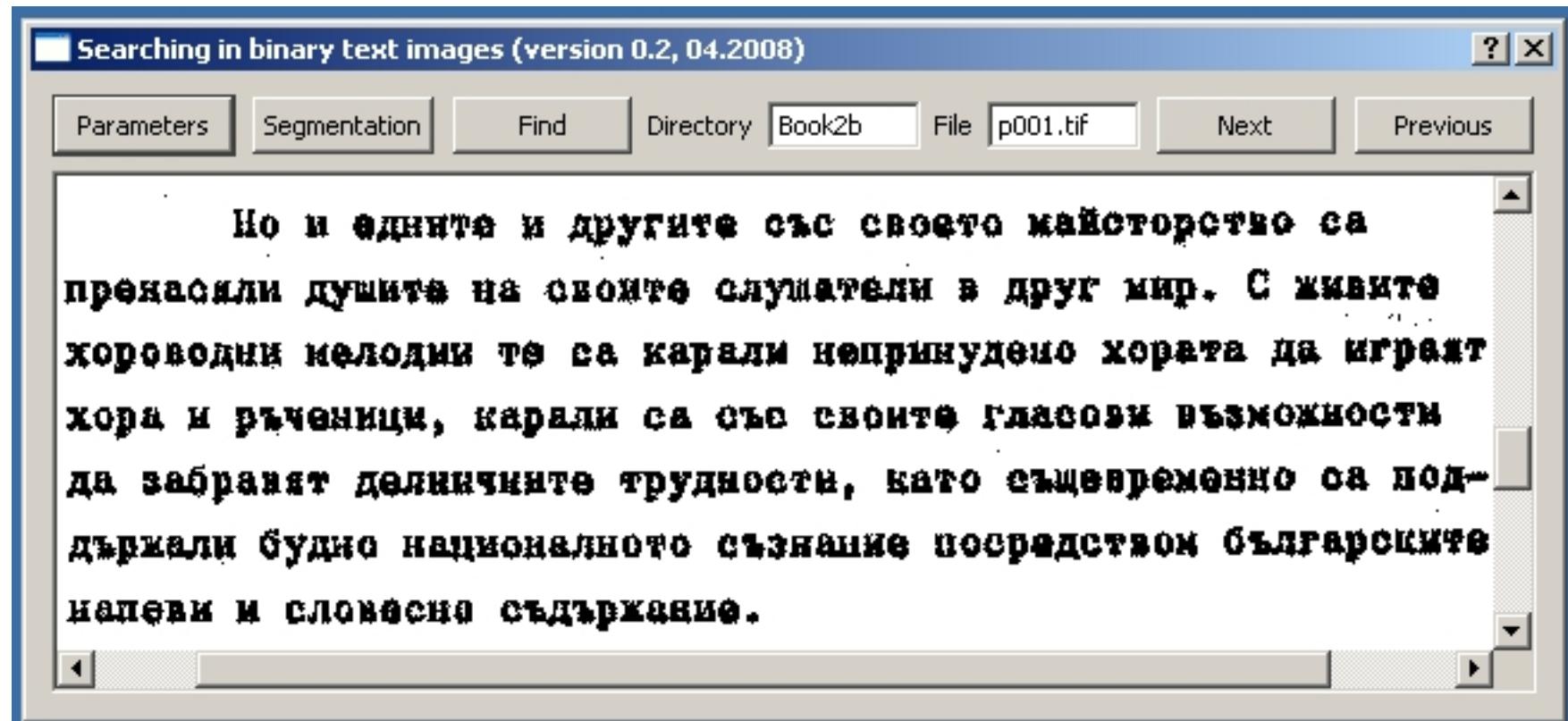
Input data of the software are collection of files representing text document.

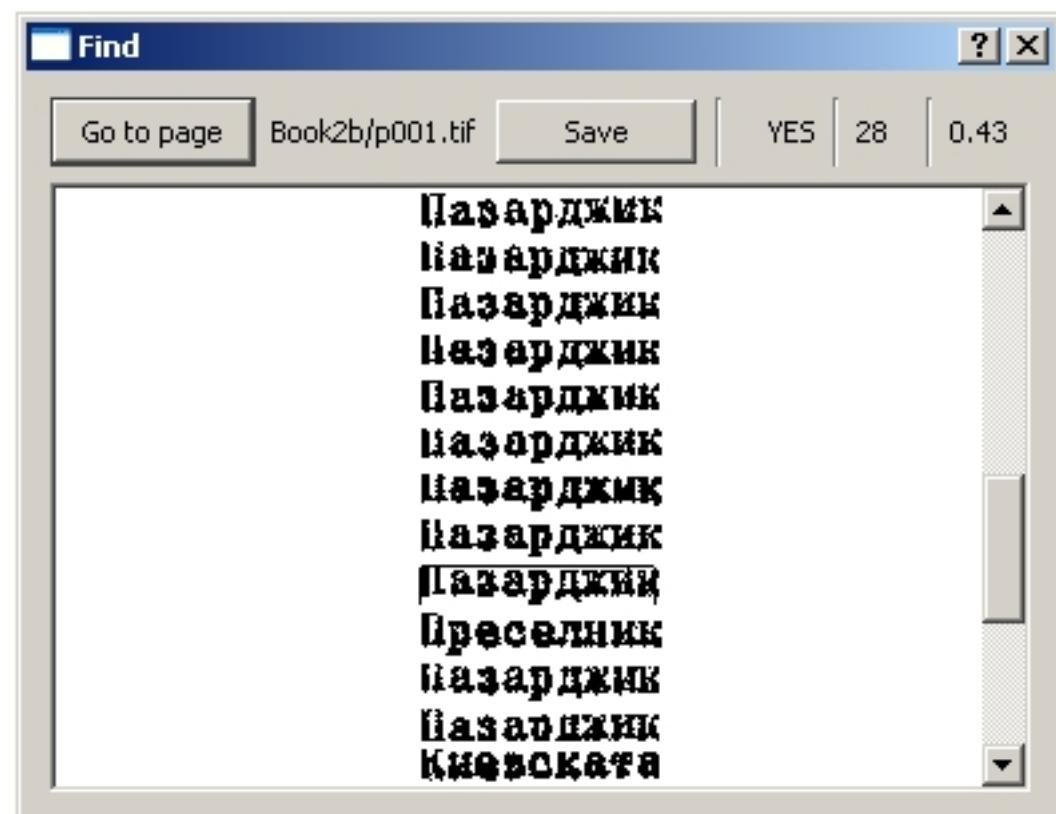
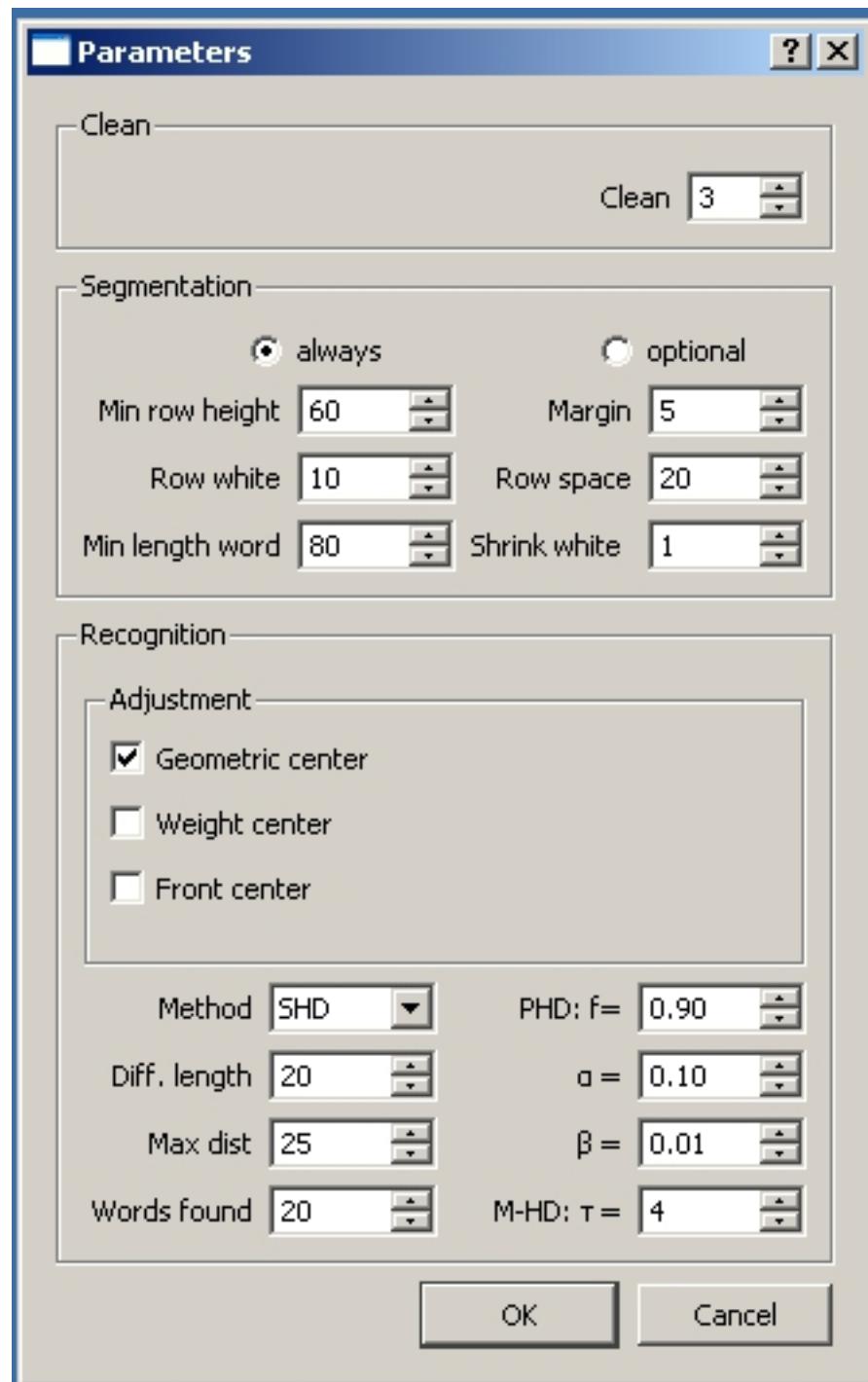
The software

The software system supports three user windows:

- Main
- Parameters
- Found.

Main window displays one page of the document – the current page. The current directory containing document image files and the current file name are given on the top of Main window. It is possible to go forward and backward through the document pages.





Segmentation

We use horizontal projection for line extraction. If the lines are horizontal straight lines or have small slopes, the histogram has near zero values between lines.

To segment the words in a line, we use vertical projection – a histogram obtained by counting the number of black pixels in each vertical scan at a given horizontal position. If the words are well separated, the histogram should have areas of zero values between words.

As a result of the segmentation, every word is associated with a word image – minimal rectangular frame that contains the corresponding word. So we consider any word image as a rectangle, which consists of white and black pixels. The black pixels form a set, which is used in calculating word similarities.

иа пред близки негови клиенти, които много обичали да го
слушат. Тремолирането на дясната му ръка е било неизминато

For segmentation step we use a number of parameters, which are important for successful word separation:

- **Minimal row height**: This parameter helps us to avoid creating (due to noise) rows with small height;
- **Margin**: This parameter allows us to process only a part of the page.
- **Row white**: When the value at a point in row histogram is less than the value of this parameter, we suppose that this point belongs to the white space between the words.
- **Row space**: Minimum value of the white space between words.
- **Minimum word length**: The system does not segment words with length less than the value of this parameter. **Shrink white**: This parameter concerns additional step conducted when we have already separated words, and words are framed. At this step we try to shrink the height of rectangles if the points of histograms have values less than this parameter.

Searching

After segmentation of a page, we must choose a pattern word image – this is a word, which we want to find in the document pages.

This step starts when the user push the button `Find` in `Main` window. It activates the process of inspection all pages to measure similarities of segmented words and the pattern word.

We can see a part of the retrieval data in `Find` window. Pushing `GoTo` button, the page containing the marked word downloads in `Main` window.

The program code is written in C++ with help of Qt – a cross-platform application development framework (open source from Trolltech).

Experiments

Bulgarian typewritten document (about 1940), 335 pages, tif (2400×3200), 1 BPP

I РАЗДЕЛ

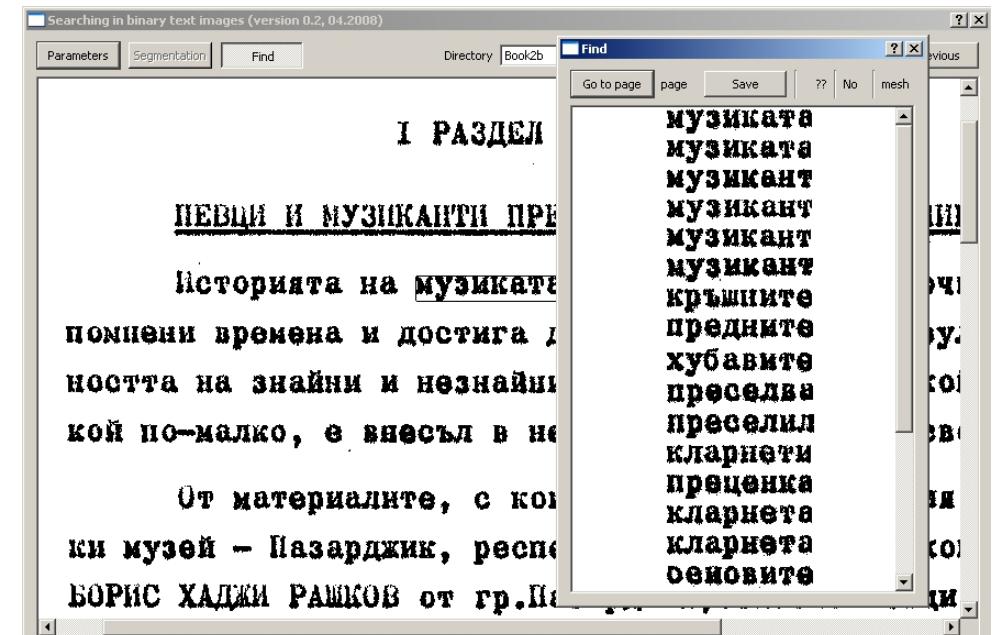
ПЕВЦИ И МУЗИКАНТИ ПРЕДИ И СЛЕД ОСВОБОЖДЕНИЕТО

Историята на музиката в гр.Пазарджик започва от незапомнени времена и достига до наши дни, като резултат от дейността на знайни и незнайни труженици, които, кой повече, кой по-малко, е внесъл в нейната съкровищница своя дад.

От материалите, с които разполага Окръжния исторически музей – Пазарджик, респективно сведенията, които е събрал БОРИС ХАДЖИ РАШКОВ от гр.Пазарджик, относно певци и музиканти преди и след Освобождението се установява, че битовите нужди, свързани с годежи, сватби, занаятчийско-еснафски сбирки, хора, вечеринки и пр. са били задоволявани от музиканти – професионалисти и любители.

Професионалисти били онези музиканти-инструменталисти или певци, като най-често инструменталиста е бил и певец, които са свирели и пеели срещу възнаграждение, а любители – онези, които със своето пеене и свирене са радвали душите и сърцата на хората по сборове, хорища и др., без да получават възнаграждение.

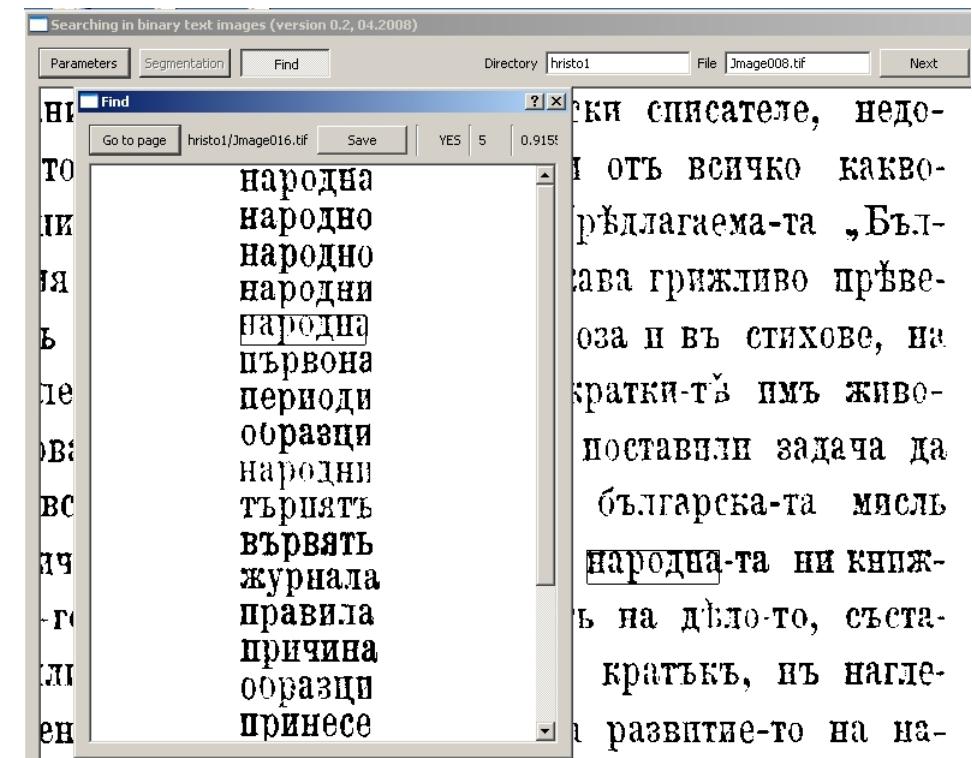
Но и едните и другите със своето маисторство са пренасяли душите на своите слушатели в друг мир. С живите хороводни мелодии те са карали непринудено хората да играят хора и ръченици, карали са със своите гласови възможности да забравят делничните трудности, като същевременно са поддържали будно националното съзнание посредством българските напеви и словесно съдържание.



Bulgarian book, Christomatia (1884), 1000 pages, tif (2300×3800), 8 BPP

Повече-то отъ ранни-тѣ му стихотворения сѫ любовни пѣсни, по подражание на грѣцки-тѣ, и не прѣдставляватъ литературина стойност; стихотворения-та му въ „Смѣсна Китка“ при всичко, че повечето сѫ слаби подражания на руски-тѣ, нѣ свидѣтелствоватъ вече за поетическо-то дарование на г. Славейкова: най-добри-тѣ му стихотворения сѫ обнародвани-тѣ пѣ-послѣ въ „Читалиште,“ отъ която „Не пѣй ми се,“ „Жестокостъ-та ми се сломи“ и „Тогасъ понѣ“ джакатъ съ истински лирпзъмъ и заслужено привлѣкоха вниманието на читателе-тѣ. Славейковъ, който е вѣнчанъ въ бѣлгарски езикъ, прѣвъ доказа гѣвкостъ-та му въ поезия-та. Като се числи между първи-тѣ борци по черковниятъ вѣпросъ, той захвашта въ сѫщото врѣме почтенно място въ редъ-тѣ на малко-то ни добри литератори.

Велико влияние е упражнила възъ пробуждането духъ-тѣ камъ свободата на независимостта у бѣлгарски народъ доста обширна-та литературна дѣятелност на *Георгий Саса Раковски* (род. въ Котелъ 1818, умр. въ Букурешть 1868 г.). Въ личността и въ дѣла-та на Раковски се отрази най-нагледно тогавашното състояние на умове-тѣ, нужди-тѣ, стремления-та и идеали-тѣ на народъ-тѣ ни. Тако-речи единичъкъ дѣнецъ по онова врѣме, той писувѣ, работї всичко. Той искаше да обгърне въ своята широка дѣялност всички-тѣ нужди на народъ-тѣ ни, да удовлетвори всички-тѣ национални купнѣния, да осѫществи най-смѣтни-тѣ и въжделени мечти. Той възвѣздаде съ фанатически вѣсторгъ минюло-то и приготви бѫдѫщите-то. Бѣше въ сѫщото врѣме поетъ, историкъ, етнографъ, публицистъ, агитаторъ и хайдутинъ. Нито на единъ бѣлгарски дѣятелъ животъ-тѣ не е билъ напълненъ съ толкова неутолима и разнообразна дѣятелност и напъстренъ съ толкова бѣди, приключения и странности. Той се бѣше училъ въ Атина, Парижъ, Цариградъ и въ Русия. Знаеше руский, срѣбский, румънский, турский, грѣцкий, староелинскай, французскай, арабский и дори отъ части санскрит-

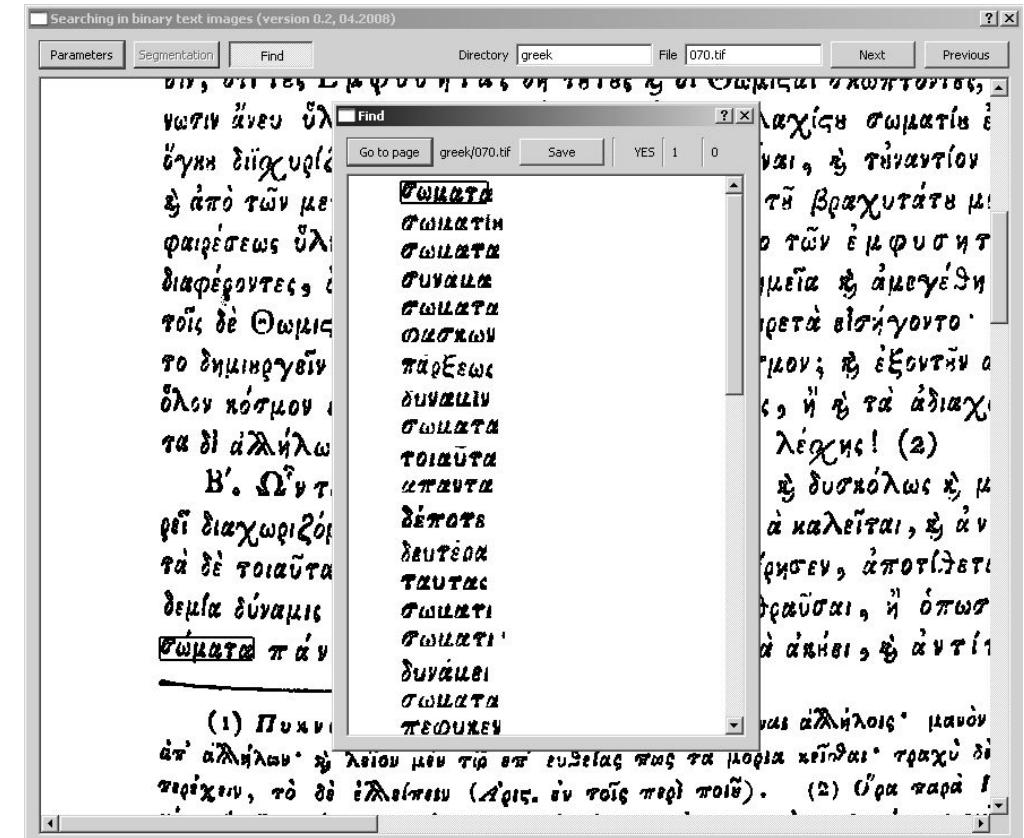


Old Greek text (approx. in the third century BC), 50 pages, jpg (1077×1416), 8 BPP

τοιαῦται κατὰ σχῆμα πάντη εἰσὶν ἄτρεπται, οἷς μ
Σῶμα δὲ σύνθετον ἐν τῇ φύσει ψδὲν τοιάτον· ὑπ
δλίψει, τοιῇ ὁπωσδήπειραι, καὶ λίθαντος οἱ σερρ
τητι διαφέρων ἀδάμας, μηδενὸς ὅλως ἔξαιρυμένη,
μοιρον, ὡς δέδεικται. Ων δὲ σωμάτων τὰ μέρη
και, καὶ ῥᾶσα διαχωρίζεσθαι πέφυκεν ἀπολαὸς ταῦτ
μέλι, ἄργιλλος, κτο· ὅσῳ δὲ ἡττονι καὶ ἀτονωτέρῳ
πλοκῆς ἀποσείχει, τοτήτῳ καὶ ἀπαλώτερῳ, ἐς ἕκρ
ψδενὶ γάρ ἐσιν ἐντυχεῖν, οὗ τὰ μέρη μὴ ὁπωσδή
γγύμενα.

Γ'. Τὸ σκληρὸν σῶμα ὑπὸ κάφιστε καὶ ἀσθενεῖ
πεψυκός, εὔθραυστον ἀκύει· τοιαῦτα χάλυψ
τὰ κεράμεια σκεύη· Τέτων τὰ εὐρεὰ μέρη ἐχότων
ἐλήλοις, διὸ καὶ δῆξε τῆς ἀμοιβαίας παρῆς ἀφίσαται

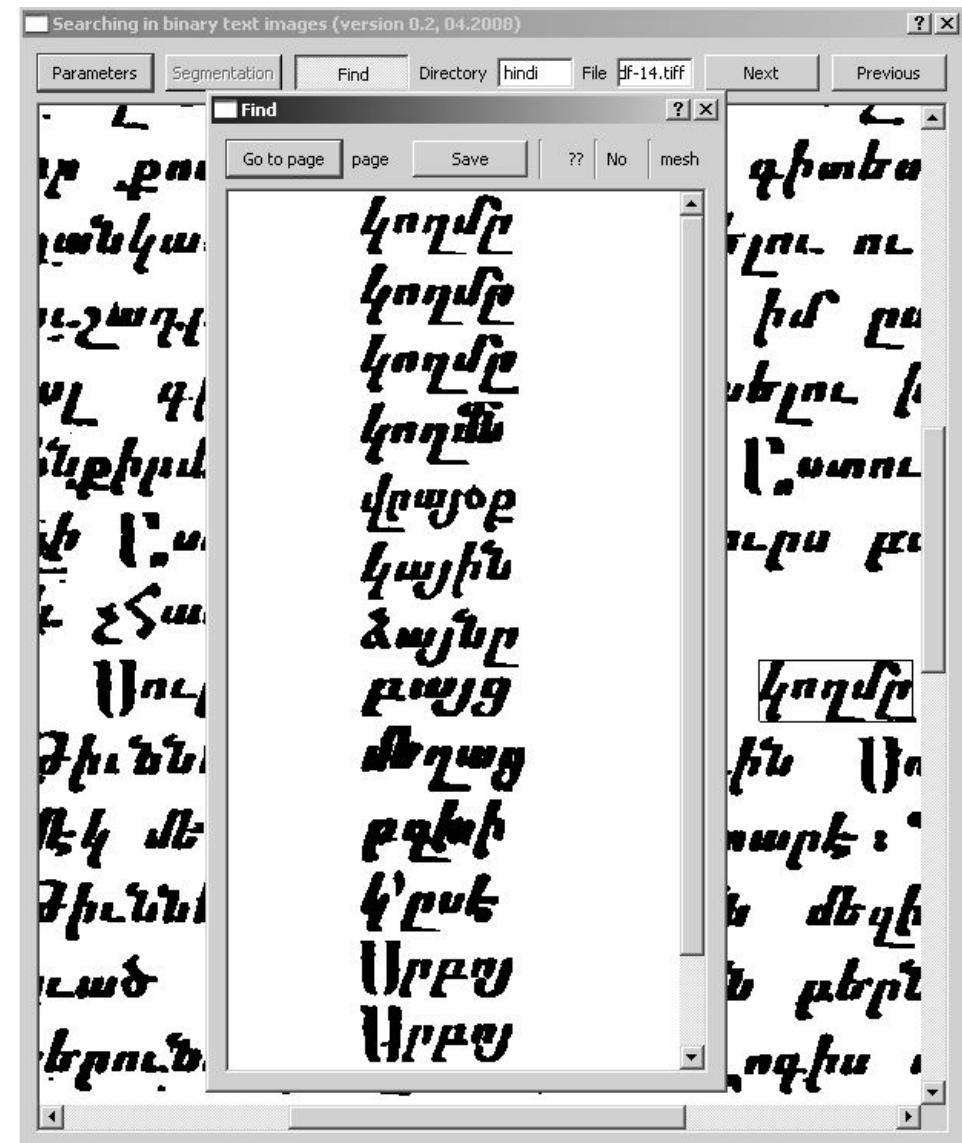
Δ'. Τὸ ἐκ πολῶν οἵουεὶ λεπίδων πάνυ λεπτῶν ἀσάμενον, εὕσχισον σῶμα προσείρηται· τότε τὰ μόρια εὐρύότερον προσκεκόληται ἀλλήλοις, ἢ λεπίς ἄριτρη τὰ τοιαῦτα σώματα εἰς λεπίδας ἀναλυόμενα· αναλέμνεται λίθοι οἱ ἐξ Ἰβηρίας, καὶ Καππαδοκίας,



Armenian book (1858), 178 pages, tiff (2800×5000), 1 BPP

կառակութեան սգին ամէն աստիճանի մարդոց սրբ-
տին մէջ տարածուեցաւ, և սորվեցուց անոնց մէկ-
զմէկ ատել ու մէկվամէկէ գարշիլ. մինչև անգամ
մէկ իանութի մէջ գործող արհեստաւորը սորվեցաւ
Նզովք տալ իր քովի խանութին մէջ բանող դրա-
ցիին, ան պատճառով որ անիկայ Հոգւցն Արբոյ
բղխումը իրեն համաձայն չդաւանիր. ոչ մէկը և
ոչ մէկալը հասկընալով թէ ինչ կ'ըսեն, կամ ինչ
բան կ'ուզեն հաստատել:

Ուստի այսպիսի պնտեղի վիճաբանութիւնները
պատճառ տուին որ երբոր մարդիկ Հոգւցն Արբոյ
վրայօք խորհին, գրեթէ միայն աս մէկ նիւթիս
ուղղեն միտքերնին, այսինքն թէ՝ Հոգւցն բղխու-
մը միայն Հօրմէն է, կամ Հօրմէն ու Արդիէն ։
Ամէնն ալ կը դաւանին թէ Հոգին Այուրբ՝ Եր-
րորդութեան մէկ անձն է. բայց ո՞վ կրնայ ըսել
թէ անիկայ ինչ ներգործութիւն կ'ընէ մարդուս
Հոգւցն փրկութեանը համար, կամ ինչ է իր
մասնաւոր պաշտօնը մարդս երկինքը բարձրացընե-
լու համար։ Եհան աս մէծ և ամենահարկաւոր
նպատակիս համար է որ Երրորդութեան վարդա-
պետութիւնը յայտնուած է։ Հայրը Խրկեց Ար-
դին աշխարհը փրկելու։ Այնուու որ Եստուած
անանկ սիրեց աշխարհը մինչև որ իր միածին Ար-



Text in Spanish (1901), 30 + 57 pages, gif (1400×2500), 4 BPP

Alonso, Rogelio M., Cartilla histyrico-descriptiva del tūrmino municipal de Macuriges. Habana:

Impr. La Propagandista, 1901, HOLLIS Catalog, Harvard University

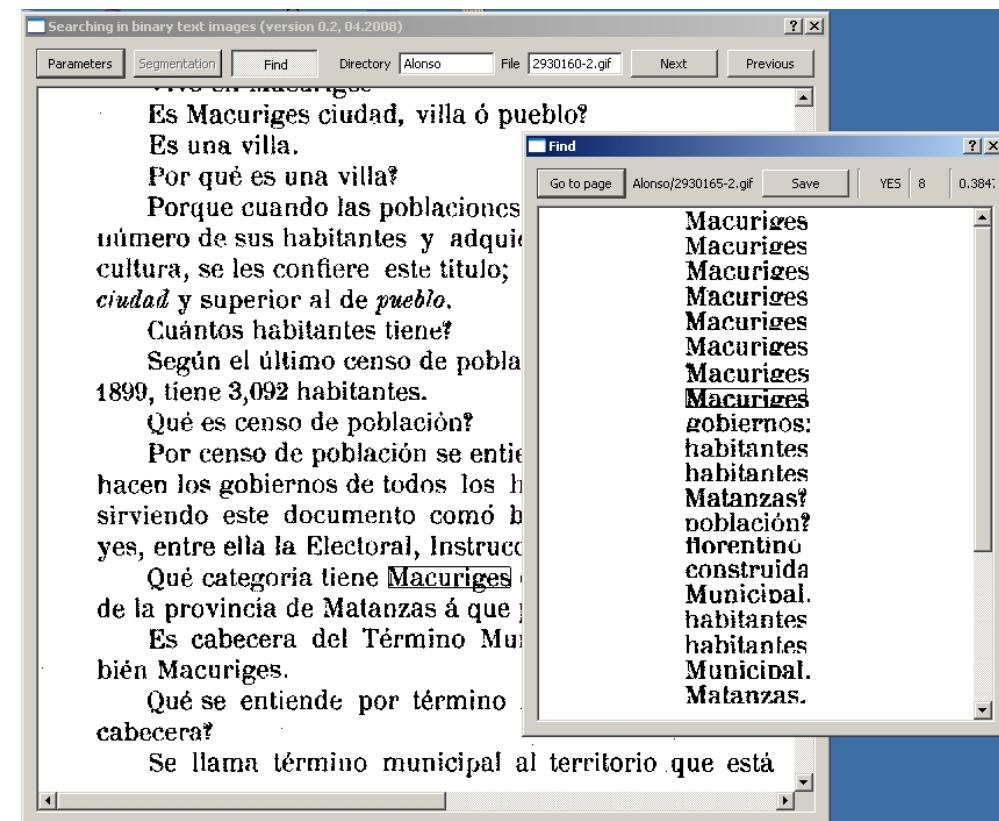
No señor; todas las fincas azucareras tienen sus chuchos que conectan con las líneas del ferrocarril y hay adcmás caminos reales, trasversales y vecinales, estos en estado natural. (1)

Qué entiende V. por caminos reales, trasversales y vecinales?

Caminos reales, son los caminos abiertos por el gobierno Español desde los tiempos primeros de la colonización de la Isla de Cuba y tienen de ancho 24 varas; caminos trasversales son los que solo tienen de ancho 12 varas y vecinales los pasos permitidos por los propietarios de fincas, para acortar distancias de un lugar á otro y salvar lo mal que pudieran estar los caminos por el fango, las piedras ó la yerba.

Cuántos ingenios para la fabricación de azúcar tiene en la actualidad el Término todo?

Los siguientes: «Santa Filomena» en el barrio de Navajas propiedad del Sr. Leandro Soler, «Elizalde» del Sr. Alberto Broch en el Ciego y «Santa Catalina» del señor Enrique Heedigg en el mismo barrio; «Carmen» del Sr. Alexander en Navajas, «Socorro» del Sr. Pedro Arenal en Tramojos y «Dolores» del Sr. Francisco Rosell en Platanal, todos centrales y con magnificos aparatos.

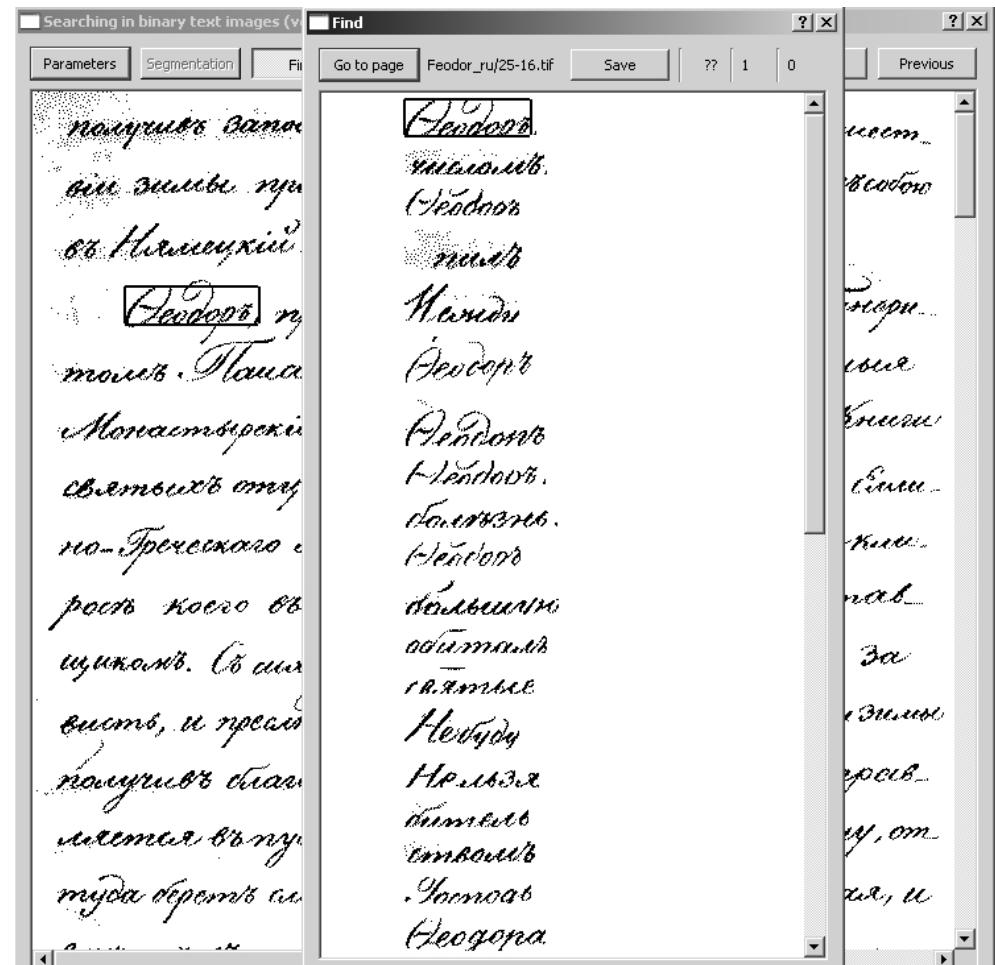


Handwritten document in Russian (1840), 44 pages, jpg (700×900), 24 BPP

Дом живоначальной Троицы, Свято-Троицкая Сергиева Лавра, Собрание славянских рукописей, 43: Житие схимонаха Феодора

Уже нога его не двигалась от бора болезни; каждая амортизаторная оледеневшая его члены. Феодоръ растворил свои кости, и наровь собственного тела, согревалъ ощущавшее тепло духовна-го друга, покрывалъ горячими лобзаниями его ѿды, освященные чистотой дѣствия, исоль-нными огнемъ Божественныхъ благодати. На рукахъ Феодора скончался великий Николай, и мощи его недаромъ прикоснулись тишине.

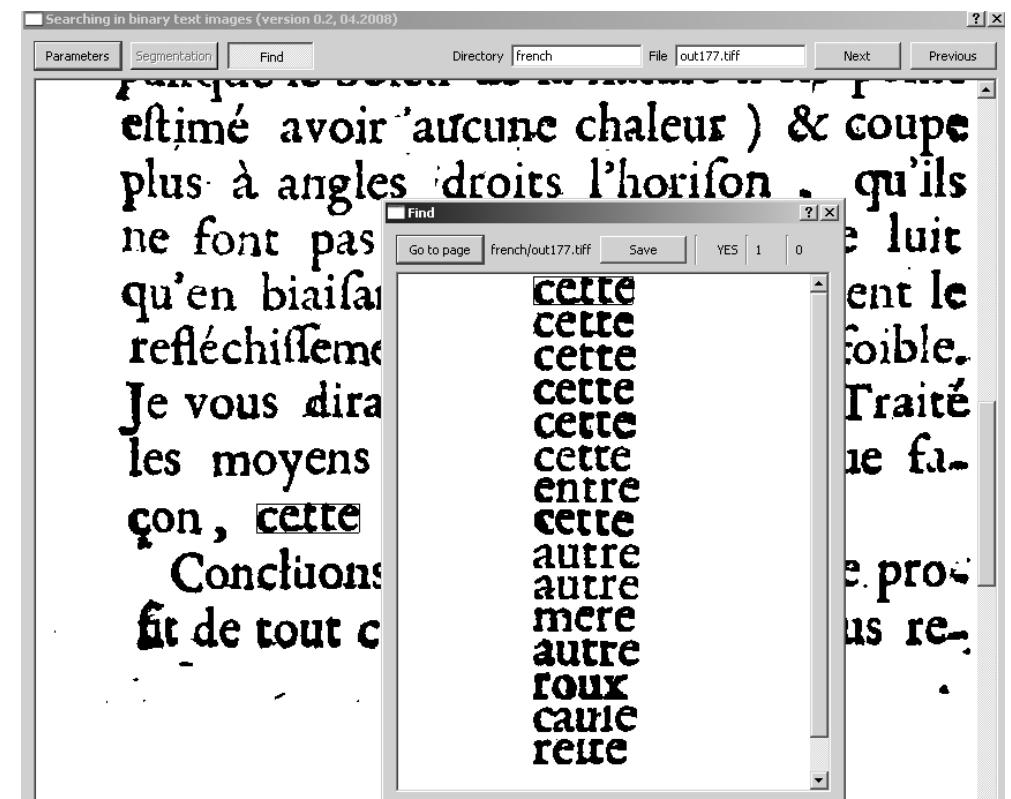
Феодоръ пребывалъ въ Никонъ до 1801^{го} года. Въ продолженіи сего времени увидѣлъ онъ кончику, высокаго чиницъ Николая, увидѣлъ кончику изнаменитаго Паниса. Преодолѣвъ сего



Text in French (1692), 388 pages, jpg (2048×3550), 8 BPP

Nicolas de Bonnefons, Ch. de Sergy, (1692), University of Gent, Digitized by Google (2007)

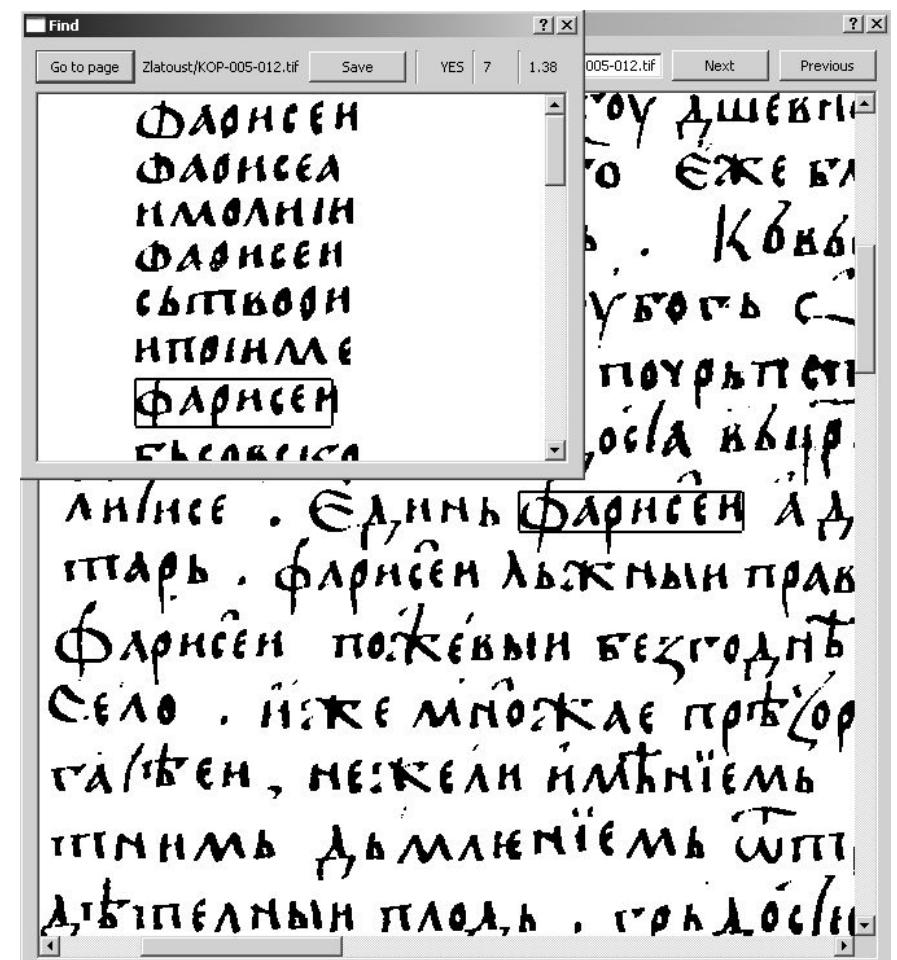
Quant à la Terre, si vous la rencontrez bonne, ce vous sera un grand avantage, & une grande épargne ; mais rarement en pourrez-vous trouver, où il n'y ait beaucoup à travailler, d'autant que telle paraîtra passablement bonne au dessus, qui étant ouverte de la profondeur d'un fer de Béche seulement, se trouvera Argileuse dessous ; ce fonds est pire aux Arbres que le Tuf, ou la Roche, à cause qu'il s'y rencontre de petites veines où les Racines peuvent s'étendre & profonder, afin de tirer la fraîcheur de plus bas, & prendre quelque nourriture ; mais l'Argileuse ou Terre franche ou rouge, fait comme un plancher qui par sa dureté & densité, ne peut être percé par aucunes Racines, & qui dans les grandes ardeurs de l'Eté, em-



Slavonic manuscript, (1574), 747 pages, jpg (1249×1890), 24 ВРР.

Дигитална Народна библиотека Србије, Ћирилски рукописи, Збирка словенских рукописа Јернеја Копитара, Зборник "Златоуст"

съмръти и мъхъ. яко даме надъющесе
богдемъ на се, пъ паба въскръшающаго
мъртвые. иже ѿполнкыє съмръти ии
збави па иизбавиъ. панже оу повахъмы
яко иеще иизбавиъ. Утъ оу бо ръци ми
къзпосищисе яко ѿ добрыи свои. Съмъ
споже благодъть исповѣдовали да въшомъ
утъ бо имаши єже непрѣель єси. аще
ли прїель єси, утосе хвалиши яко непрїе
мъ. не ты ба позналъ єси правдою, пъ
вътебъ благослію позна. вѣще рече ба,
пакъ же познани въвше ѿ ба. не ты ба
прїель єси добродѣтелю. пъ тебъ хсъ
пришъствїемъ прїеть. гонеко рѣ аще
и постнгну имже и постнжанъ въихъ
ѡхъ. невѣ мене и избрасте рѣ гъ. пъ
азъ и избралъ се. пъ ли зане поуѣтень єси
веле мудръствѹеши; имѣть въвни



Conclusion

We show the possibilities of preliminary version of the program.

Experiments with 7 different languages from various epochs give us the certainty in practical benefit of our approach.

It is quite universal and does not required any specific features of the concrete language.

Word searching can be applied to any collection of scanned documents, immediately after the graphic files have been created by the scanner device.

In spite of relatively low efficiency of the Hausdorff type methods (the searching process takes a lot of time) high level personal computers will be able to solve the problem in reasonable time.

We think also that the accuracy of retrieval is quite sufficient for practice.

There are some directions to improve the software.

- Increasing the efficiency and speeding up the search.
- Searching with a part of word as a pattern.
- Character segmentation of a page (or pages) and composing pattern word from well separated letters.
- Feedback – making second search for the same word with a different pattern word.
 - The user can choose this word among correct words found in the first search.
 - Produce a new pattern as an average of all or part of the correct words.
- Automatic or semi automatic choice of parameters based on image information.

Thank you for the attention.