

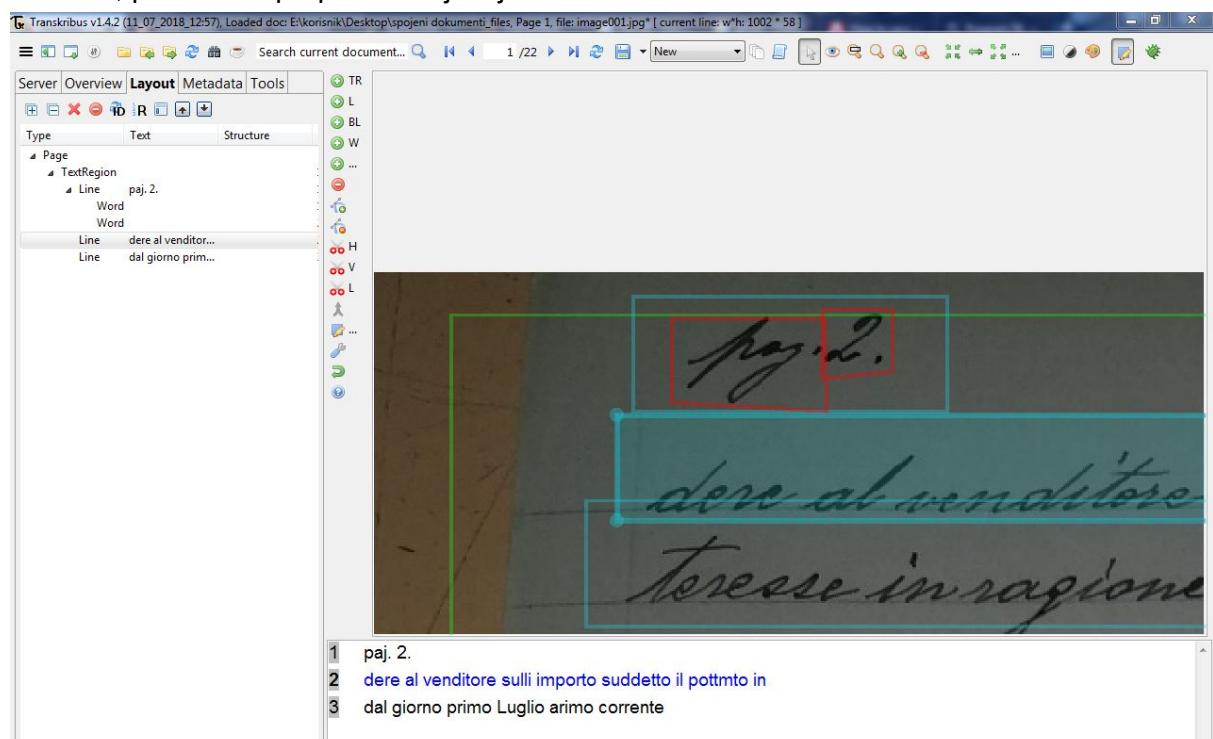
Doktorand Josip Mihaljević prisustvovao je 18. listopada 2018. radionici *Transkribus: optičko prepoznavanje znakova rukopisnih tekstova*. U prvoj dijelu radionice, koji se održao u Rektoratu sveučilišta u Zagrebu, Damir Boras, Vlatka Lemić i Hrvoje Stančić govorili su o problemima digitalizacije te optičkoga prepoznavanja znakova i međunarodnim projektima povezanim s digitalizacijom i upravljanjem digitalnim dokumentima poput [ICARUS-a](#).



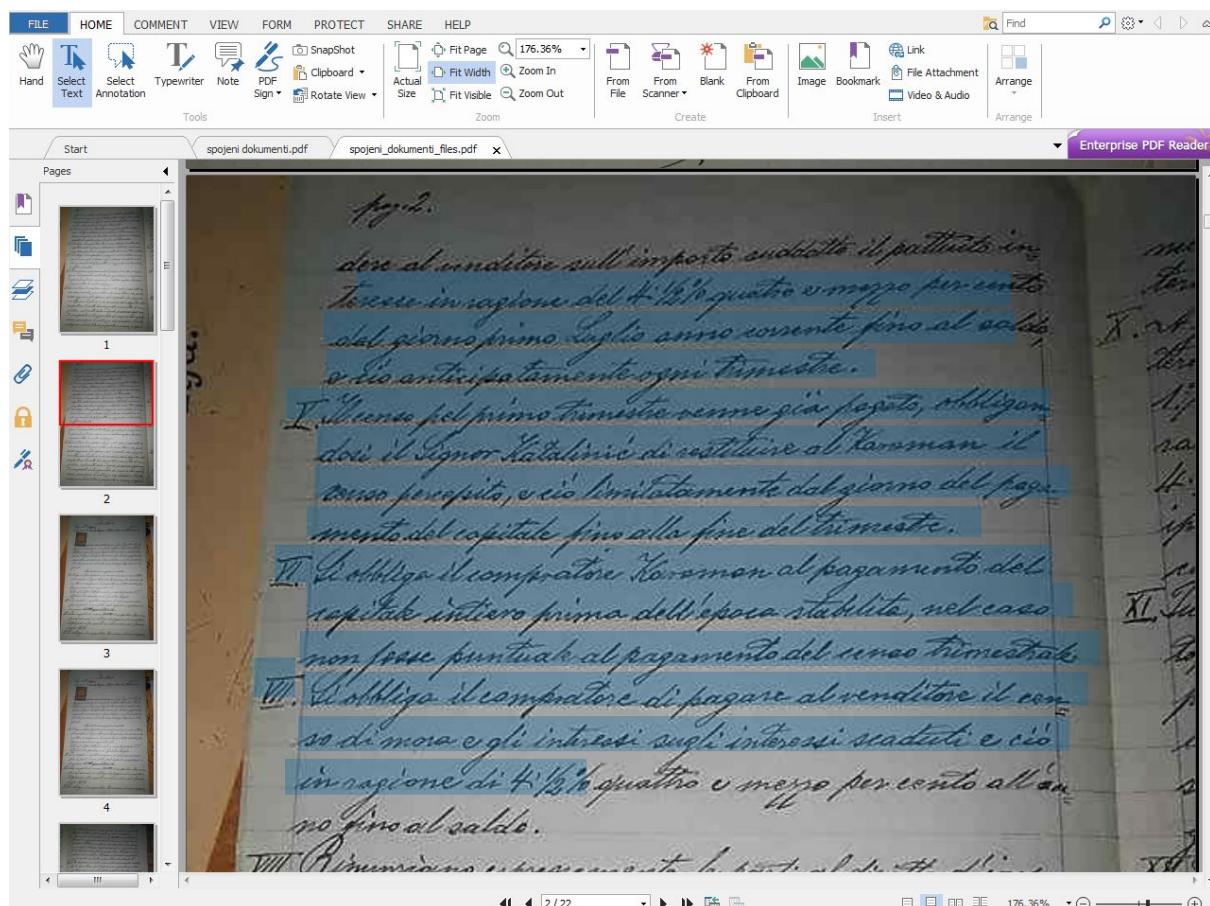
Unutar projekta [READ](#), koji je financiran iz fondova Europske unije, razvio se program Transkribus koji je izrađen za optičko prepoznavanje rukopisnih znakova. Program je nastao kako bi se pronašao način za automatsko optičko prepoznavanje znakova povijesnih dokumenata. Program pokušava automatski prepoznati slova rukopisa na temelju unesenoga skupa podatka kojima se koristi za trening te dodatno omogućuje pretraživanje sadržaja digitaliziranih rukopisa. Podržava i stare rukopise poput hebrejskoga, cirilice i arapskoga. Program se može besplatno preuzeti te trenutačno ima oko 15 000 korisnika. U programu se pohranjuju slike i prijepisi koji se koriste za trening prepoznavanja znakova. Program može točniji prepoznati znakove ako ima više primjera određenoga rukopisa koji prepoznaće. Uglavnom se program trenira da prepozna rukopise određenih osoba. Moguće je stvoriti i modul koji sadržava rukopise više osoba, ali sigurnije je imati odvojene module za prepoznavanje znakova svakoga autora. Preciznost je programa za svaki znak oko 90 posto ako program sadržava dovoljno velik uzorak stranica s obrađenim rukopisom. Pogreške u prepoznavanju znakova mogu se ručno ispraviti. Pretraživati se može s pomoću ključnih riječi. Kao primjer uporabe toga programa spominje se obrada 75 000 zapisu Jeremyja Bentham u Londonu. Digitalizirano je 95 000 stranica za koje je napravljen ručni prijepis koji se koristio za treniranje programa da prepozna ostale znakove skeniranih Benthamovih rukopisa. Za označavanje određenih znakova na stranicama koristio se TEI XML. Bentham je s vremenom oslijepio, pa je njegov rukopis postao manje čitljiv te ga je bilo teže prepoznati. Zbog toga su napravljena dva odvojena modula za trening njegovih rukopisa. Napravljena je [stranica](#) na kojoj se može mrežno pretraživati rukopise. Upute za rad s programom napisane su na [wiki stranicama](#). Spominje se i tehnički proces rada u programu te problemi poput izbjedjelih i umrljanih znakova u rukopisu. Često se pojavljuju i tekstovi napisani u dva stupca te tekstovi u tablicama koji se teže mogu automatski obraditi. Kao

algoritam prepoznavanja koristi se Neural Network, koji pamti sekvene podataka za prepoznavanje znakova. *Keyword spotting* tehnologija je koja omogućuje automatski pronađenje određenih riječi u slikama rukopisa. Koristi se u programu za pronalaženje određenih rukopisa koji sadržavaju tu riječ. Inačica program koji je dostupan na službenoj stranici zove se [TranskribusX](#). Predstavljen je i program za digitalizaciju na mobitelu DocScan, koji u kombinaciji s programom TentScan izravno prebacuje slike na Transkribus.

U drugom dijelu radionice radio se u programu Transkribus. Za korištenje programa prvo je potrebna registracija s pomoću Googleova računa. Svaki polaznik trebao je učitati svoje slike rukopisa te označiti na njima linije teksta za koje se trebalo napraviti prijepis. Kasnije se utvrdilo da postoji opcija za automatsko označavanje linija u tekstu te da se mogu upotrijebiti postojeći istrenirani moduli za prepoznavanje teksta kako bi se automatski napravilo optičko prepoznavanje teksta. Nažalost, nije bilo mogućnosti prepoznavanja za hrvatskih rukopisnih tekstova, pa točnost prepoznavanja nije bila dobra.



Da bi se napravio vlastiti modul za optičko prepoznavanje znakova određenog rukopisa, korisnik mora za najmanje 50 skeniranih odkumenat imati precizno označene riječi s prijepisom. Nakon toga može se kontaktirati proizvođače programa da stvore modul za prepoznavanje znakova na temelju postojećega uzorka. Obrađene rukopise može se izvesti u različitim formatima poput datoteka Word, .txt, .xml te PDF-a. U PDF dokumentu na slikama rukopisa moguće je kopirati tekst.



Program također ima u sebi integriran program Abby Finereader za optičko prepoznavanje tiskanih znakova. Taj program podržava hrvatski jezik te ima visoku točnost. Program nije besplatan. Postoji, međutim, mrežna inačica koja ima ograničenja u vezi s korištenjem. Unutar Transkribusa može se koristiti besplatno bez ograničenja.

Transkribus v1.4.2 (11_07_2018_12:57), Loaded doc: slikice.ID: 89405, Page 3, file: Capture.jpg [Image Meta Info: [Resolution:96.0, w*h: 671 * 685]] [current word: w*h: 55 * 16]

Server Overview Layout Metadata Tools

Type Text Structure

Page Printspace

TextRegion Poštovani.

TextRegion ovaj tim studenata pohađa kolegij Društveno korisno učenje u informacijskim znanostima u zimskom semestru ak. god. 2017/18 na Odsjeku za Informacijske znanosti Filozofskog fakulteta Sveučilišta u Zagrebu. Studenti viših godina već 10 godina sudjeluju na projektima društveno korisnog učenja u kojima neprofitnim udružama i javnim institucijama stavljuju na raspolaganje svoje informacijske i informatičke vještine i znanja tijekom cijelog semestra.

Kolegij Društveno korisno učenje kroz literaturu, uključenost u projekt u neprofitnoj udruži ili javnoj instituciji ili kroz vođenje dnevnika iada o rasprodavanju na projektu omogućava studentima viših godina studija Informacijskih znanosti ozbiljnu i konkretnu implementaciju znanja i vještina koja su stekli tijekom studija pružajući partnerskoj instituciji besplatnu realizaciju potrebnog informatičkog / informacijskog rješenja s ciljem lakšeg snalaženja na prvom poslu po završetku studija.

Studenti posjeduju različita znanja koja se mogu primjeniti na različitim zadatcima počevši od informacijskog opisovanja različitih dobnih skupina, priprema publikacija za tisk, grafičkog dizajna, izrada baza podataka, rada s animacijama sve do pouduke osoblja partnerske institucije za rad u određenim programima i sa potrebni softverom.

Kao dio obveza kolegija, svaki student će u dogovoru s predstavnikom partnerske institucije sudjelovati u projektu koji će se cijeli semestar odvijati u instituciji u studentskom timu. Tamski projekti čini najvažniji dio kolegija i od svakog studenta se očekuje 40 sati rada na projektu, koji mora biti zavrsen do kraja semestra.

Informacijska oprema koja bi studentima bila na raspolaganju je poželjna, ali ne i nuzna. Konačan proizvod studenata je besplatan, kao i njihov rad i savjeti. Studenti svoj projekt mogu raditi i kod kuće, uz dogovor s Vama i mentorom na fakultetu. Svi podaci iz Vaše institucije koji se stavljuju na raspolaganje studentima smatraju se povjerljivima, i mogu se dati na uvid samo studentima koji sudjeluju na projektu i studentskom centru.

Molimo Vas da pomognete u realizaciji ovog procesa učenja i rasta, omogućavajući našim studentima izvođenje projekta za Vasu ustanovu.

1 ovaj tim studenata pohađa kolegij Društveno korisno učenje u informacijskim znanostima
2 u zimskom semestru ak. god. 2017/18. na Odsjeku za Informacijske znanosti Filozofskog
3 fakulteta Sveučilišta u Zagrebu. Studenti viših godina već 10 godina sudjeluju na
4 projektima društveno korisnog učenja u kojima neprofitnim udružama i javnim
5 institucijama stavljuju na raspolaganje svoje informacijske i informatičke vještine i
6 znanja tijekom cijelog semestra.

Polaznici su isprobali i Android aplikaciju [DocScan](#) za optičko prepoznavanje znakova, koja zahtijeva korištenje uređaja *ScanTent*. *ScanTent* funkcioniра kao mali šator u koji se položi dokument ili knjiga za skeniranje. Na vrhu šatora nalazi se podloga s rupom za postavljanje mobilnoga uređaja. Program *DocScan* automatski radi digitalizaciju za svaku novo prepoznatu stranicu. Može primjetiti pokrete ruke te tek uključuje skeniranje kad korisnik mirno položi papir. Skenirani dokumenti mogu se učitati za obradu na stranicama Transkribusa.

