

SVEUČILIŠTE U ZAGREBU  
FILOZOFSKI FAKULTET  
ODSJEK ZA INFORMACIJSKE ZNANOSTI  
Ak. god. 2009./ 2010.

Ante Kranjčević

## **APLIKACIJA ZA RAČUNANJE N-GRAMA**

Diplomski rad

Mentor: dr. sc. Kristina Vučković

Zagreb, 2010.

# Sadržaj

Sadržaj .....	1
1. Uvod .....	2
2. Teorijski dio .....	4
2.1. n-grami znakova.....	4
2.2. Markovljev model .....	5
2.3. Frekvenčijska analiza .....	5
3. <i>bigtri</i> aplikacija .....	7
3.1. Izrada aplikacije.....	7
3.2. Unos i statistika teksta.....	7
3.2.1. Metoda unosa i obrade teksta.....	7
3.2.2. Stuktura podataka u bazi .....	10
3.2.3. Unos obrađenih podataka u bazu .....	12
3.2.4. Vanjsko spremanje rezultata obrade teksta .....	13
3.2.5. Statistički prikaz rezultata obrade teksta .....	14
3.2.6. Vanjsko spremanje statističkih podataka .....	16
3.3. Statistika teksta .....	16
3.3.1. Statistika <i>top30</i> .....	17
3.3.2. <i>Top30</i> vanjsko spremanje rezultata .....	18
4. Usporedba s drugim aplikacijama .....	19
4.1. Wolfram Alpha .....	19
4.1.1. Unos i interpretacija .....	19
4.1.2. Prikaz rezultata obrade teksta.....	20
4.1.3. Vanjska obrada rezultata obrade teksta.....	24
4.2. WordCreator .....	26
4.2.1. Unos i interpretacija .....	26
4.2.2. Prikaz rezultata obrade teksta.....	27
4.2.3. Vanjska obrada rezultata obrade teksta.....	28
5. Ocjena razlika .....	29
5.1. Unos teksta .....	29
5.2. Prikaz podataka obrade teksta .....	29
5.3. Brzina obrade teksta .....	30
5.4. Vanjska obrada rezultata .....	30
5.5. Konačna ocjena .....	31
6. Rezultati pojavljivanja bigrama i trigramu .....	33
7. Zaključak .....	36
8. Popis literature .....	37

# 1. Uvod

Nemoguće je u današnje doba zamisliti obradu teksta bez pomoći računala. Velika količina teksta već se nalazi zapisana u digitalnom obliku što ga čini pogodnjim za računalnu obradu, za razliku od teksta koji se nalazi u pismenom obliku kojeg je prethodno potrebno digitalizirati. Računalna, kao i ručna obrada teksta, vrši se na mnogim razinama ovisno o korisničkim potrebama. Ovaj rad opisuje mrežnu aplikaciju nazvanu *bigtri* čija zadaća je obrada teksta na način da iz njega izdvaja bigrame<sup>1</sup>, trigrame<sup>2</sup>, bigrame riječi<sup>3</sup> i trigrame riječi<sup>4</sup>.

Bigrami, trigrami, bigrami riječi i trigrami riječi lingvističke su jedinice koje su sastavni dio jezika kojim se služimo. U informacijskom riječniku za njih koristimo skupnu riječ n-grami. Na osnovu n-grama jezici se uspoređuju na lingvističkoj bazi, te se izrađuje tablica učestalosti pojavljivanja istih lingvističkih elemenata u svakom pojedinom jeziku.

Misao vodila, te svrha izrade ovoga rada bila je nepostojanje takvoga alata na hrvatskom jeziku prema mojim saznanjima, te pojednostavljinjanje rada korisnicima koji se koriste ovakvom analizom jezika. Aplikacija je postavljena na mrežni server tako da joj je moguće pristupiti u bilo kojem trenutku, dostupna je svima bez restrikcija, što povećava količinu teksta u bazi. Veća količina spremlijenog teksta daje nam preciznije izlazne podatke.

Korisnici u radu s aplikacijom nisu ograničeni obrađivanjem samo hrvatskog jezika, nego je moguć izbor 38 svjetskih jezika koji su prikazani u tablici 1.

Ovaj je rad podijeljen u dva dijela. U prvom dijelu detaljno su pojašnjene metode izrade aplikacije, te svaki dio aplikacije koji korisnik koristi u radu. U drugom dijelu rada napravljena je usporedba sa sličnim postojećim alatima koji se mogu pronaći besplatno na internetu. Prilikom ocjenjivanja i uspoređivanja aplikacija u obzir su uzeti funkcionalnost, brzina obrade unešenog teksta, izgled obrađenih podataka, te mogućnost obrade izlaznih podataka.

---

<sup>1</sup> *bigram* – niz od dva uzastopna znaka

<sup>2</sup> *trigram* – niz od tri uzastopna znaka

<sup>3</sup> *bigram riječi* – niz od dvije uzastopne riječi

<sup>4</sup> *trigram riječi* – niz od tri uzastopne riječi

<b>Jezik</b>				
Albanski	Filipinski	Japanski	Poljski	Švedski
Armenijski	Finski	Kineski	Portugalski	Talijanski
Bjeloruski	Francuski	Latvijski	Rumunjski	Tibetanski
Bosanski	Grčki	Litvanski	Ruski	Turski
Danski	Hrvatski	Makedonski	Slovački	Ukrajinski
Engleski	Indonezijski	Mongolski	Slovenski	Vijetnamski
Esperanto	Irski	Norveški	Srpski	
Estonski	Islandska	Njemački	Španjolski	

**Tablica 1 – Popis jezika u bazi**

Podatke dobivene analizom jezika (pojavljivanje pojedinih bigrama, trigrama, bigrama riječi, trigrama riječi) moguće je koristiti kao referentne podatke za daljnje istraživanje jezika.

## 2. Teorijski dio

### 2.1. *n-grami znakova*

N-grami znakova su pod-nizovi  $n$  predmeta iz danog niza znakova. Predmeti u ovom slučaju mogu biti fonemi<sup>5</sup>, slogovi riječi ili riječi. N-gram veličine 1 naziva se *unigram*, n-gram veličine 2 naziva se *bigram*, n-gram veličine 3 naziva se *trigram*, dok se n-grami veličine 4 ili više jednostavno nazivaju *n-grami*. Modeli n-grama su tipovi modela koji se koriste za predviđanje sljedeće stavke u svakom slijedu, te se koriste u različitim područjima statističkih obrada prirodnih jezika. N-grami modela modeliraju sekvence, posebice prirodnih jezika, pri tome koristeći statistička svojstva n-grama. Primjeri korištenja statističkih svojstava n-grama su još u prepoznavanju govora, obradi fonema kao i obrađivanja sekvenci fonema.

Statistička vjerojatnost sljedećeg n-grama bila je ideja za eksperiment Claudea Shannona<sup>6</sup> u radu o teoriji informacija 1948. godine. Shannon je postavio pitanje, kako s obzirom na redoslijed slova možemo izračunati vjerojatnost sljedećeg slova u riječi. Promatrajući podatke došao je do odgovora da može izvesti distribuciju vjerojatnosti za sljedeće slovo koje ima zbroj vjerojatnosti svih sljedećih slova 1.0. Npr. za sekvencu "for ex" sljedeći je n-gram po njegovom izračunu imao vrijednost  $a = 0,4$ ,  $b = 0,00001$ ,  $c = 0, \dots$  i tako sva ostala slova abecede dok zbroj vjerojatnosti njihovih pojavljivanja nije 1.0. Ovakvi modeli n-grama često su kritizirani u literaturi jer im nedostaje eksplisitni prikaz dalekometne ovisnosti<sup>7</sup>. Najžešći kritičar ovoga modela je *Noam Chomsky*<sup>8</sup>, čija kritika je proizašla iz razloga što modeli n-grama koji su spomenuti ranije u tekstu imaju raspon ovisnosti ( $n-1$ ), te se kao takvi ne mogu koristiti kao rješenje za neograničene ovisnosti koje stvaraju prirodni jezici. Iz tog razloga n-gram modeli nisu imali veći utjecaj na razvoj lingvističke teorije nego se koriste samo u praktičnim aplikacijama za obradu jezika, glasovno prepoznavanje jezika i generiranje slučajnih nizova fonema (Manning, Schütze, 1999.).

<sup>5</sup> *Fonem* – najmanja jedinica govora koja je bitna za značenje

<sup>6</sup> *Claude Shannon* (1916. – 2001.) – Američki matematičar i ing. elektrotehnike, naziva se ocem informacijske teorije

<sup>7</sup> *Dalekometna ovisnost* – fenomen koji se javlja u analizama prostornih i vremenskih nizova podataka, a odnosi se na stopu raspadanja statističke ovisnosti

<sup>8</sup> *Noam Chomsky* (1928.) – Američki filozof, lingvist i politički aktivist, otac generativne gramatike

## **2.2. *Markovljev model***

Markovljev lanac, nazvan po Andreyu Markovu<sup>9</sup>, predstavlja niz stanja sustava koji se koristi za statističko modeliranje. Sustav koji je opisao Markov takav je da u svakom trenutku može prijeći u neko novo stanje ili može ostati u istome stanju što ovisi samo o trenutnom stanju u kojemu se sustav nalazi.

Markovljevi lanci imaju širok spektar uporabe u statistici, književnosti, informacijskoj znanosti. Primjer Markovljevog lanca je sljedeći: ako pratimo prehrambene navike osobe koja jede samo grejp, sir i grašak i prati sljedeća pravila da jede samo jednom dnevno i ako je jučer jeo sir, neće ga jesti danas, dok će grašak i grejp jesti uz istu vjerojatnost. Ako je jučer jeo grejp, danas će jesti grejp uz vjerojatnost  $1/10$ , sir  $4/10$  i grašak  $5/10$ . Konačno, ako je jučer jeo grašak, neće ga jesti danas, no jesti će grejp uz vjerojatnost  $4/10$  ili sir  $6/10$ . Prehrambena navika iz primjera može se prikazati Markovljevim lancem, jer odluka što će danas jesti, ne ovisi o tome što je osoba jela prije nekoliko dana, nego jedino ovisi o tome što je jučer jela (Rabiner, Juang, 1986.).

Ovakvu proceduru koristio je i Claude Shannon, prilikom predviđanja mogućnost pojavljivanja n-grama, a to je upravo Markovljev lanac. Shannon je eksperimentom pokazao da i je bez poznavanja cijelog sustava moguće pribaviti mnogo statističkih pravilnosti istog tog sustava.

## **2.3. *Frekvencijska analiza***

Frekvencijska analiza koristi se kod prikazivanja učestalosti znakova, slova, bigrama, trigrama i ostalih komponenti teksta u jeziku. Određena pravilnost se pojavljuje u nabrojanim slučajevima, te je pogodna za analizu jer postoji karakteristična distribucija bigrama ili trigrama karakteristična za svaki pojedini jezik. Kao primjer navest ćemo engleski jezik u kojemu je slovo „e“ veoma često, dok se s druge strane slovo „x“ ne pojavljuje tako često u jeziku. Što se bigrama tiče, u engleskom jeziku, bigrami 'st', 'ng', 'th' veoma su česti, za razliku od bigrama 'nz' ili 'qj' (Wikipedia).

---

<sup>9</sup> Andrey Markov (1856. – 1922.) – Ruski matematičar najpoznatiji po radu na teoriji stohastičkih procesa koji su poslije nazvani markovljev lanac

Analize pokazuju kako frekvencija pojedinih znakova, odnosno bigrama i bigrama riječi, varira ovisno o autoru, kao i o temi o kojoj autor piše. Kao primjer možemo uzeti tekst u kojemu se piše o „x-zračenju“, te u tom tekstu možemo očekivati veliku frekvenciju znaka 'x', odnosno bigrama koji u sebi sadrži znak 'x'. Znakovi, odnosno slova, te, bigrami, trigrami, kao i bigrami, odnosno trigrami riječi mogu se koristiti kao dokaz autentičnosti nekoga teksta, te dokazati da li je pojedino djelo napisao autor za kojeg se to tvrdi. Preciznu prosječnu frekvenciju pojavljivanja bigrama, tragrama, bigrama riječi i tragrama riječi možemo dobiti samo analizom velike količine reprezentativnog teksta (Manning, 1999.).

Frekvencija pojavljivanja znakova imala je velik utjecaj na dizajn pojedinih tipkovnica, tj. rasporeda tipki na njima. Kao primjer pojavljuje se takozvana *AZERTY* tipkovnica, koja je specifična po tome što su najfrekventnija slova stavljeni u prvi red tipkovnice (*AZERTY*), za razliku od tipkovnice kakvu mi koristimo kojoj je raspored slova u gornjem redu *QWERTZ* (Wikipedia).

### **3. *bigtri* aplikacija**

U sljedećem poglavlju pokazat će na koji je način ustrojena aplikacija, kako funkcioniraju pojedine komponente programskog kôda, te funkcionalnost i snalaženje u korištenju same aplikacije. Aplikacija je pojednostavljena i organizirana tako da bude što jednostavnija za uporabu, korisniku pruža maksimalnu funkcionalnost i omogućava što jednostavnije snalaženje prilikom korištenja.

#### **3.1. Izrada aplikacije**

Aplikacija *bigtri* izgrađena je u C# programskom sučelju koristeći Microsoft-ovu ASP. net tehnologiju servera kako je već spomenuto u uvodu. Koristeći programski jezik C# dobio sam jednostavnost i veliku mogućnost implementacije s ostalim komponentama koje su korištene u izradi ovog alata. Brzina obrade jezika, korištenje potencijala za rad nad bazom podataka, te velika baza korisničke podrške elementi su koji su prevagnuli u izboru programskog alata za izradu *bigtri* aplikacije.

#### **3.2. Unos i statistika teksta**

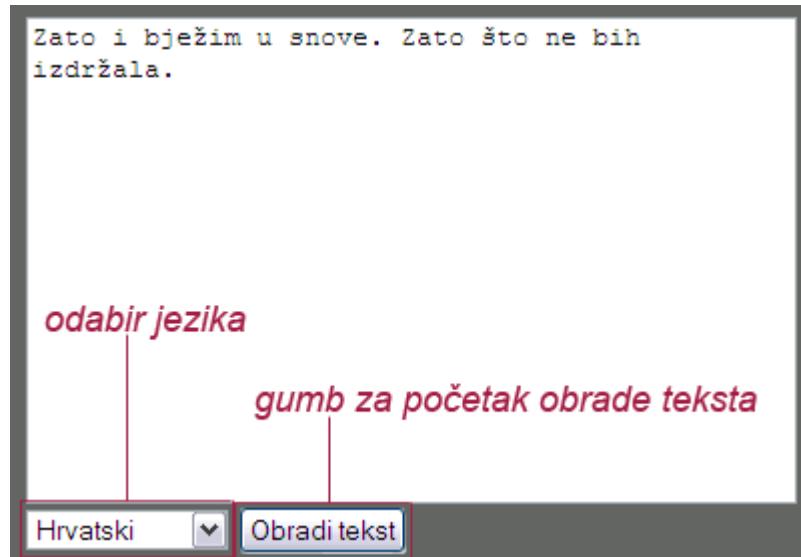
##### **3.2.1. Metoda unosa i obrade teksta**

Tekst koji se obrađuje i analizira unosi se preko komponente za unos teksta koje se nalazi na početnoj mrežnoj stranici [www.ante.sibinj.com](http://www.ante.sibinj.com). Nakon unosa teksta, bilo ručnim unosom ili *copy*<sup>10</sup> - *paste*<sup>11</sup> metodom unosa, potrebno je odabrati jezik na kojem je tekst napisan (korisnik ima na izbor 38 jezika kao što je prikazano u tablici 1), te pritisnuti gumb *Obradi tekst*, kako je prikazano na slici 1. Pritiskom na gumb aktivira se dio programskog kôda koji preko jednostavne petlje iz teksta izdvaja jedan po jedan bigram, trigram, bigram riječi, te trigram riječi. Niz znakova koji dobijemo i u koji su spremljeni bigrami, trigrami, bigrami riječi i trigrami riječi spremaju se u listu koja se koristi tijekom dalnjeg rada aplikacije (izrada statistike, spremanje u tekstualnu datoteku).

---

<sup>10</sup> *copy* – kopirati

<sup>11</sup> *paste* – zalijepiti



Slika 1 – Forma za unos teksta

Prilikom pisanja petlji za obradu teksta potrebno je bilo voditi brigu o slučajevima kada se u tekstu umjesto jednog razmaka između riječi nalazi više razmaka. Provjera razmaka veoma je bitna, jer inače ne bih dobio pravilne bigrame, trigrane, bigrame riječi i trigrane riječi nakon obrade teksta. Programska petlja prikazana je na slici 2.

U ovom dijelu koda jasno se vidi kako naredba *Replace* kroz While petlju provjerava cijeli tekst. Petlja se izvršava onoliko puta koliko je potrebno dok ne ukloni sve višestruke razmake između riječi.

```
while (victim.IndexOf(" ") >= 0)
{
    victim = victim.Replace(" ", " ");
}
victim = victim.Replace(Environment.NewLine, "");
```

Slika 2 – Petlja za provjeru razmaka

Primjer formatiranja, te pravilnog i nepravilnog računanja bigrama prikazan je u tablici 2. Kao primjer naveo sam rečenicu „**Zato i bježim.**“ U rečenici iz primjera se ispred slova 'i' nalazi dvostruki razmak, umjesto pravilnog jednostrukog razmaka.

Zato i bježim.	
Nepravilno	Pravilno
'Za'	'Za'
'at'	'at'
'to'	'to'
' '	' i'
' i'	' i '
'i '	'bj'
'bj'	'je'
'je'	'ež'
'ež'	'ži'
'ži'	'im'
'im'	'm.'
'm.'	

**Tablica 2 – Primjer pravilne i nepravilne obrade teksta**

Iz tablice se vidi da je razlika između pravilno obrađenog teksta i nepravilno obrađenog teksta u jednom bigramu. Razlika nastaje zbog jednog znaka viška, tj. zbog dodatnog razmaka koji se pojavljuje u rečenici. Jasno je da takvi slučajevi nikako ne dolaze u obzir ukoliko želimo točan broj obrađenih leksičkih jedinica.

Petlja koja obrađuje tekst uklanja sve višestruke razmake i pretvara ih u jedan razmak kako bi tekst bio pravilno obrađen. Osim dvostrukih razmaka, petlja obrađuje slučajeve kada se između susjedna dva reda teksta nalazi više od jednog praznog reda pretvarajući ih u jedan red. Ručni unos jednog ili više uzastopnih praznih redova spriječen je onemogućavanjem pritiska tipke *enter* prilikom ručnog unosa teksta u komponentu za unos teksta. Ovo ograničenje se vidi iz sljedećeg koda koji je prikazan na slici 3.

```
<asp:TextBox ID="txtText" runat="server"
    Height="239px" TextMode="MultiLine"
    Width="381px"
    onkeydown = "return (event.keyCode!=13);">
</asp:TextBox>
```

**Slika 3 – Ograničenje pritiska tipke *enter***

Slika 4 prikazuje kôd kojim biramo jezik iz padajućeg izbornika.

```
<asp:DropDownList ID="ddlJezik" runat="server"
    DataSourceID="AccessDataSource1"
    DataTextField="Jezik" DataValueField="ID">
</asp:DropDownList>
```

**Slika 4 – Izbor jezika iz padajućeg izbornika**

Kôd koji obrađuje tekst pokreće se pritiskom na gumb *Obradi tekst* kako je i navedeno ranije u poglavlju. Izgled kôda prikazan je na slici 5.

```
<asp:Button ID="btnObradi" runat="server"
    onclick="btnObradi_Click" Text="Obradi tekst"/>
```

**Slika 5 – Kôd obrade teksta**

### 3.2.2. Stuktura podataka u bazi

Obrađene podatke potrebno je spremiti u bazu podataka kako bi se mogli koristiti u dalnjem radu. Aplikacija koristi *Microsoft Access* bazu podataka za spremanje podataka. Access baza odabrana je zbog jednostavnosti rada. U ovom je slučaju baza sačinjena od jedne datoteke koja u sebi sadrži šest povezanih tablica navedenih u tablici 3.

Tablice		
Jezik	BigramChar	TrigramChar
Tekst	BigramWord	TrigramWord

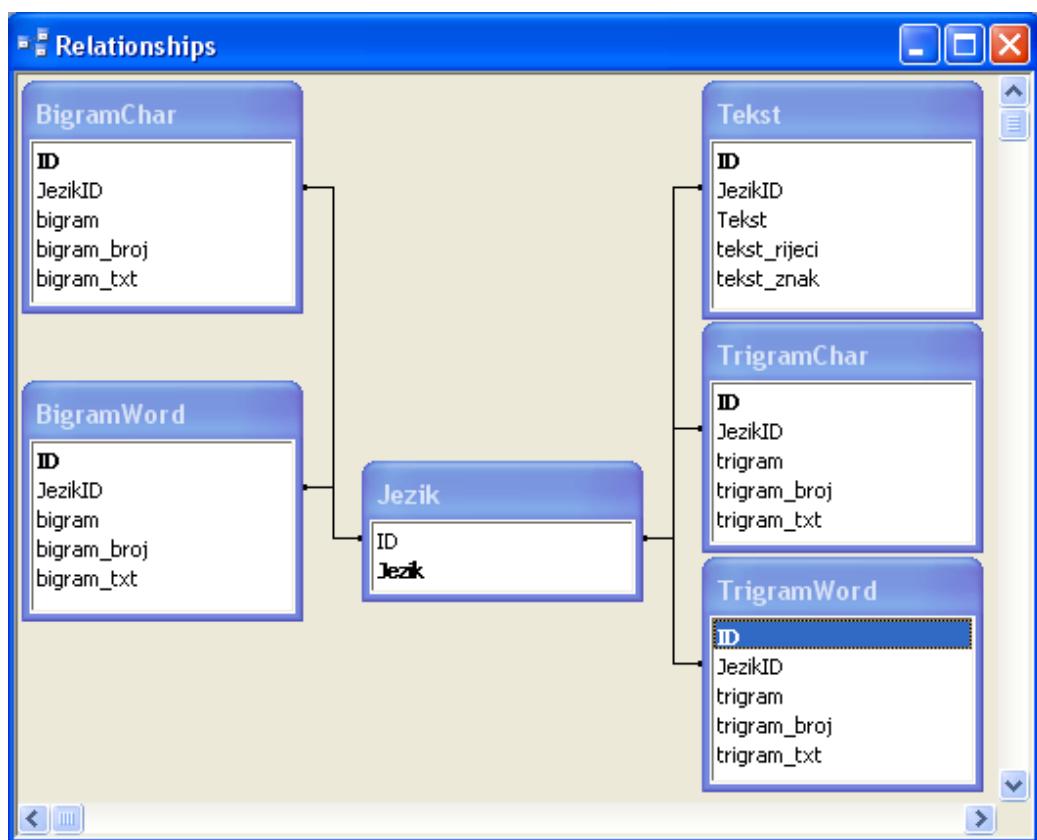
**Tablica 3 – Popis tablica u bazi podataka**

Alternativa Access-u, je *SQL*<sup>12</sup> baza podataka, no ona zahtijeva posebni server za rad s njom i iz tog razloga nije korištena u aplikaciji. Odabirom Access baze podataka niti u jednom trenutku nije narušena funkcionalnost. Jedini mogući

<sup>12</sup> *SQL* – kratica za *Structured Querry Language* – računalni jezik dizajniran za rad nad bazama podataka

nedostatak je brzina rada s bazom, jer se SQL baza podataka odlikuje većom brzinom pristupa i obrade podataka u tablicama.

Tablica *Jezik* glavna je tablica u bazi. Ona integrira sve ostale tablice preko jedinstvenog ključa. Jedinstveni ključ polje je naziva *ID* (jedinstveni identifikacijski broj svakog pojedinog jezika) koje je bročanog tipa i u kojemu se ne mogu pojavljivati dvostrukе vrijednosti. Nemogućnost pojavljivanja dvostrukе vrijednosti u polju znači da svaki jezik u tablici ima jedinstveni identifikacijski broj pomoću kojega identificiramo ostale podatke u bazi. Jezici iz ove tablice vidljivi su korisnicima u padajućem izborniku na početnoj stranici aplikacije.



Slika 6 – Integritet tablica unutar baze podataka

Tablica *Tekst* služi za spremanje svih tekstova koji su upisani u komponentu za unos teksta i nakon toga obrađeni. Na taj način, ova tablica služi kao kontrolna tablica, kako bi se u svakom trenutku mogli ponovno provjeriti rezultati dobiveni nakon obrade teksta. Ostale četiri tablice (*BigramChar*, *BigramWord*, *TrigramChar*, *TrigramWord*) služe za spremanje bigrama, trigrama, bigrama riječi i trigrama riječi pojedinačno u svaku od navedenih tablica. Logička struktura i integritet tablica unutar baze podataka prikazani su na slici 6.

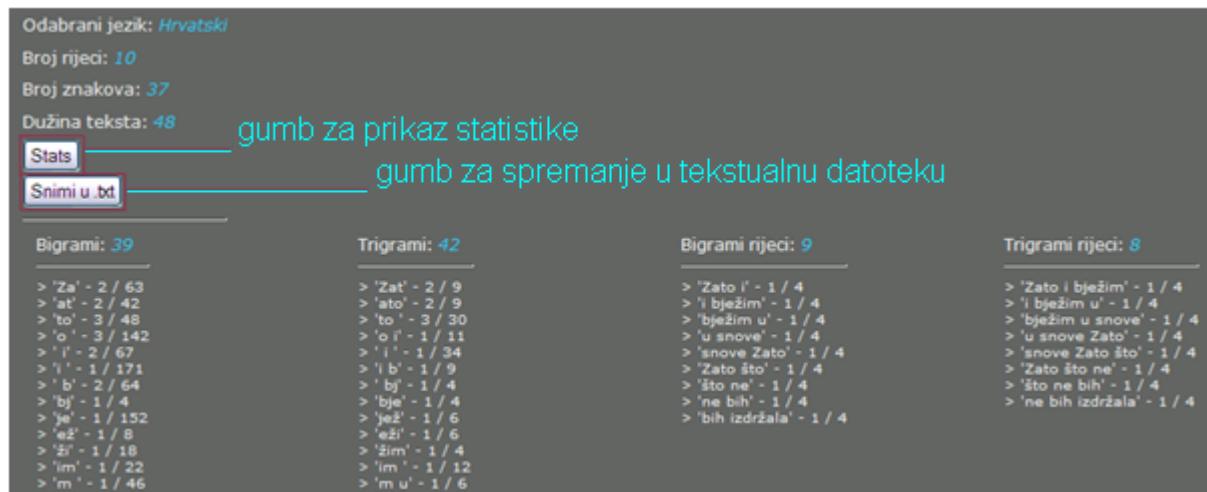
### 3.2.3. Unos obrađenih podataka u bazu

Nakon unosa i obrade podataka, sljedeća zadaća aplikacije je spremiti podatke u bazu. Spremanje započinje na način da se lista podataka koja je dobivena nakon obrade teksta, kao što je opisano u poglavlju 3.2.1., uspoređuje s podatcima koji se već nalaze u svakoj od pojedinih tablica za odabrani jezik. Ukoliko se u tablicama već pojavljuju bigrami, trigrami, bigrami riječi i trigrami riječi, koji su dobiveni i nakon obrade teksta, ne unose se u nova polja u tablice, nego se njihov broj pojavljivanja uvećava za broj pojavljivanja dobiven u obrađenom tekstu. S druge strane, ukoliko se radi o prvom pojavljivanju istih (bigrama, trigrama, bigrama riječi i trigrami riječi), unose se u tablice u nova prazna polja. Kao početni broj pojavljivanja postavlja se broj dobiven iz obrađenog teksta kako je prikazano na slici 7.

```
if(iNgramID != 0)
{
    using (OleDbCommand cmd = new OleDbCommand("UPDATE " + tableToSaveTo + " SET " + column +
    "_broj = @ngram_broj, " + column + "_txt = @ngram_txt WHERE ID = @ID", conn))
    {
        cmd.Parameters.Add(new OleDbParameter("@ngram_broj", iBigramBroj + ngram.Count));
        cmd.Parameters.Add(new OleDbParameter("@ngram_txt", iBigramTxt + BrojZnakova));
        cmd.Parameters.Add(new OleDbParameter("@ID", iNgramID));
        cmd.ExecuteNonQuery();
    }
    ngram.CountBefore = iBigramBroj;
}
else
{
    using (OleDbCommand cmd = new OleDbCommand("insert into " + tableToSaveTo + "(JezikID, " + column +
    ", " + column + "_broj, " + column + "_txt) values (@JezikID, @ngram, @ngram_broj, @ngram_txt)", conn))
    {
        cmd.Parameters.Add(new OleDbParameter("@JezikID", jezikID));
        cmd.Parameters.Add(new OleDbParameter("@ngram", ngram.Znakovi));
        cmd.Parameters.Add(new OleDbParameter("@ngram_broj", ngram.Count));
        cmd.Parameters.Add(new OleDbParameter("@ngram_txt", BrojZnakova));
        int i = cmd.ExecuteNonQuery();
    }
    ngram.CountBefore = 0;
}
```

Slika 7 – Provjera postojanja leksičkih jedinica u bazi podataka

Rezultati obrade unešenog teksta vidljivi su ispod komponente za unos teksta na početnoj stranici aplikacije. Izgled rezultata obrade teksta prikazan je na slici 8.



Slika 8 – Rezultat obrade teksta

Rezultati obrade grupirani su na način da u prvom stupcu imamo ispisani ukupan broj bigrama obrađenog teksta, te svaki pojedini bigram jedan ispod drugoga. Pored svakog od dobivenih bigrama nalazi se broj pojavljivanja istoga u obrađenom tekstu, te broj pojavljivanja ukoliko se bigram od prije nalazi spremlijen u bazi. Kao primjer naveo sam bigram 'Za'<sup>13</sup> koji se u obrađenom tekstu pojavljuje 2 puta, dok je u bazu spremlijen već 63 puta. Brojevi pojavljivanja bigrama međusobno su odvojeni znakom „/“. Na isti način nabrojani su i ostali elementi teksta i to redom trigrami, bigrami riječi i trigrami riječi.

### 3.2.4. Vanjsko spremanje rezultata obrade teksta

Osim unosa obrađenih podataka u bazu i njihova uspoređivanja, moguće je podatke spremiti i u vanjsku datoteku pritiskom na gumb *Snimi u .txt* koji se nalazi na vrhu stranice. Pritiskom na gumb podatci se na računalo spremaju u datoteku početnog naziva *export.txt*. Kôd je prikazan na slici 9.

```
<asp:Button ID="btnSnimi" runat="server"
Text="Snimi u .txt"
onclick="btnSnimi_Click" />
```

Slika 9 – Spremanje statistike u vanjsku datoteku

<sup>13</sup> Postoji razlika između 'Za' i 'za' bigrama, te ih aplikacija upisuje u posebna polja u tablicu

Datoteka je jednostavna tekstualna datoteka koju je moguće otvoriti u najjednostavnijem računalnom programu za obradu teksta koji dolazi instaliran u sklopu operacijskog sustava na računalima. Mogućnost vanjskog spremanja podataka dodana je jer velik broj korisnika želi nakon obrade teksta dobiti obrađene podatke u tekstualnom formatu kako bi ih mogli koristiti na osobnom računalu. Izgled datoteke nakon otvaranja u tekstualnom *editoru*<sup>14</sup> prikazan je na slici 10.

```

export.txt - Notepad
File Edit Format View Help
Odabrani jezik: Hrvatski
Broj riječi: 10
Broj znakovna: 3/
Bigrami: 39; Bigrami riječi: 9; Trigrami: 42; Trigrami riječi: 8
'zā' - 2 / 59; 'zato i' - 1 / 2; 'zat' - 2 / 5; 'Zato i bježim' - 1 / 2
'at' - 2 / 38; 'i bježim' - 1 / 2; 'ato' - 2 / 5; 'i bježim u' - 1 / 2
'to' - 3 / 42; 'bježim u' - 1 / 2; 'to' - 3 / 24; 'bježim u snove' - 1 / 2
'u' - 3 / 136; 'u snove' - 1 / 2; 'o i' - 1 / 9; 'u snove Zato' - 1 / 2
'i' - 2 / 63; 'snove Zato' - 1 / 2; 'i i' - 1 / 32; 'snove Zato Što' - 1 / 2
'i' - 1 / 169; 'zato što' - 1 / 2; 'i b' - 1 / 7; 'Zato što ne' - 1 / 2
'b' - 2 / 60; 'što ne' - 1 / 2; 'bj' - 1 / 2; 'što ne bih' - 1 / 2
'bj' - 1 / 2; 'ne bih' - 1 / 2; 'bje' - 1 / 2; 'ne bih izdržala' - 1 / 2
'jc' - 1 / 150; 'bih izdržala' - 1 / 2; 'jež' - 1 / 4;
'ež' - 1 / 6; 'eži' - 1 / 4;
'ži' - 1 / 16; 'žim' - 1 / 2;
'im' - 1 / 20; 'im' - 1 / 10;
'm' - 1 / 44; 'm u' - 1 / 4;
'u' - 1 / 57; 'u' - 1 / 31;
'u' - 1 / 100; 'u s' - 1 / 9;
's' - 1 / 120; 'sn' - 1 / 2;

```

Slika 10 – Rezultati obrade u tekstualnom obliku

### 3.2.5. Statistički prikaz rezultata obrade teksta

Nakon obrade vezane uz unos bigrama, trigrama, bigrama riječi i trigrama riječi, moguće je prikazati statistiku i usporedbu istih s bigramima, trigramima, bigramima riječi, te trigramima riječi koji se nalaze spremljeni u bazi. Usporedba je moguća ne samo za odabrani jezik, nego i za preostale jezike iz baze (njih 38).

Do statističkih podataka dolazimo pritiskom gumba *Stats* koji se pojavljuje na početnoj stranici aplikacije nakon obrade teksta (slika 4). Kôd kojim otvaramo statistiku prikazan je na slici 11.

```

<asp:Button ID="btnBigramStats"
    runat="server"
    onclick="btnBigramStats_Click" Text="Stats" />

```

Slika 11 – Spremanje statistike u vanjsku datoteku

<sup>14</sup> *editor* – aplikacija za uređivanje određenih tipova podataka

Nakon otvaranja stranice pojavljuje se tablica s lijeve strane prozora u kojoj su sadržani usporedni podaci za preostale jezike iz baze. Osim rezultata obade teksta, u tablici se nalaze i nabrojani bigrami, trigrami, bigrami riječi i trigrami riječi koji su u bazu ranije unešeni i obrađeni, ali za druge jezike iz baze. Pored svake leksičke jedinice nalazi se njezin broj pojavljivanja, kao i statistički podatak koji govori koja je frekvencija njezina pojavljivanja u pojedinom jeziku. Izgled stranice koja sadrži tablicu podataka prikazan je na slici 12.

bigtri stats																																																																													
Odabran jezik: Hrvatski																																																																													
Broj riječi: 10																																																																													
Broj znakova: 37																																																																													
<a href="#">Snimi u .txt</a>																																																																													
gumb za spremanje u tekstualnu datoteku																																																																													
<table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 5%;">Jezik:</td> <td>Albanski</td> </tr> <tr> <td>Jezik:</td> <td>Armenijski</td> </tr> <tr> <td>Jezik:</td> <td>Byeloruski</td> </tr> <tr> <td>Jezik:</td> <td>Bosanski</td> </tr> <tr> <td>Jezik:</td> <td>Bugarski</td> </tr> <tr> <td>Jezik:</td> <td>Danski</td> </tr> <tr> <td>Jezik:</td> <td>Engleski</td> </tr> <tr> <td>i - 1</td> <td>0,04167e b - 10,03333</td> </tr> <tr> <td>b - 1</td> <td>0,04167</td> </tr> <tr> <td>im - 2</td> <td>0,08333</td> </tr> <tr> <td>e - 2</td> <td>0,08333</td> </tr> <tr> <td>Jezik:</td> <td>Esperanto</td> </tr> <tr> <td>Jezik:</td> <td>Estonski</td> </tr> <tr> <td>Jezik:</td> <td>Filipinski</td> </tr> <tr> <td>Jezik:</td> <td>Finski</td> </tr> <tr> <td>Jezik:</td> <td>Francuski</td> </tr> <tr> <td>Jezik:</td> <td>Grčki</td> </tr> <tr> <td>Jezik:</td> <td>Indonezijski</td> </tr> <tr> <td>Jezik:</td> <td>Irski</td> </tr> <tr> <td>Jezik:</td> <td>Islandska</td> </tr> <tr> <td>Jezik:</td> <td>Japanski</td> </tr> </table>		Jezik:	Albanski	Jezik:	Armenijski	Jezik:	Byeloruski	Jezik:	Bosanski	Jezik:	Bugarski	Jezik:	Danski	Jezik:	Engleski	i - 1	0,04167e b - 10,03333	b - 1	0,04167	im - 2	0,08333	e - 2	0,08333	Jezik:	Esperanto	Jezik:	Estonski	Jezik:	Filipinski	Jezik:	Finski	Jezik:	Francuski	Jezik:	Grčki	Jezik:	Indonezijski	Jezik:	Irski	Jezik:	Islandska	Jezik:	Japanski																																		
Jezik:	Albanski																																																																												
Jezik:	Armenijski																																																																												
Jezik:	Byeloruski																																																																												
Jezik:	Bosanski																																																																												
Jezik:	Bugarski																																																																												
Jezik:	Danski																																																																												
Jezik:	Engleski																																																																												
i - 1	0,04167e b - 10,03333																																																																												
b - 1	0,04167																																																																												
im - 2	0,08333																																																																												
e - 2	0,08333																																																																												
Jezik:	Esperanto																																																																												
Jezik:	Estonski																																																																												
Jezik:	Filipinski																																																																												
Jezik:	Finski																																																																												
Jezik:	Francuski																																																																												
Jezik:	Grčki																																																																												
Jezik:	Indonezijski																																																																												
Jezik:	Irski																																																																												
Jezik:	Islandska																																																																												
Jezik:	Japanski																																																																												
<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 25%;">Bigrami: 39</th> <th style="width: 25%;">Trigrami: 42</th> <th style="width: 25%;">Bigrami riječi: 9</th> <th style="width: 25%;">Trigrami riječi: 8</th> </tr> </thead> <tbody> <tr> <td>- Za - 61 <b>0,11531</b></td><td>- Zat - 7 <b>0,00632</b></td><td>- Zato i - <b>0,00143</b></td><td>- Zato i bježim - <b>0,00263</b></td></tr> <tr> <td>- at - 40 <b>0,07561</b></td><td>- ato - 7 <b>0,00632</b></td><td>- 3 <b>0,00143</b></td><td>- i bježim - <b>0,00263</b></td></tr> <tr> <td>- to - 45 <b>0,08507</b></td><td>- to - 27 <b>0,02439</b></td><td>- i bježim <b>0,00143</b></td><td>u - 3 <b>0,00263</b></td></tr> <tr> <td>- o - 139 <b>0,26726</b></td><td>- o i - 10 <b>0,00903</b></td><td>- 3 <b>0,00143</b></td><td>- bježim u - <b>0,00263</b></td></tr> <tr> <td>- i - 65 <b>0,12287</b></td><td>- i - 33 <b>0,02981</b></td><td>- bježim u <b>0,00143</b></td><td>slove - 3 <b>0,00263</b></td></tr> <tr> <td>- b - 62 <b>0,11720</b></td><td>- ib - 8 <b>0,00723</b></td><td>- 3 <b>0,00143</b></td><td>- u slove <b>0,00263</b></td></tr> <tr> <td>- bj - 3 <b>0,00567</b></td><td>- bj - 3 <b>0,00271</b></td><td>- 3 <b>0,00143</b></td><td>Zato - 3 <b>0,00263</b></td></tr> <tr> <td>- je - 151 <b>0,28544</b></td><td>- bje - 3 <b>0,00271</b></td><td>- anove <b>0,00143</b></td><td>- anove <b>0,00263</b></td></tr> <tr> <td>- ež - 7 <b>0,01323</b></td><td>- jež - 5 <b>0,00452</b></td><td>Zato - 3 <b>0,00143</b></td><td>Zato što - <b>0,00263</b></td></tr> <tr> <td>- ži - 17 <b>0,03214</b></td><td>- eži - 5 <b>0,00452</b></td><td>- Zato što <b>0,00143</b></td><td>3 <b>0,00263</b></td></tr> <tr> <td>- im - 21 <b>0,03970</b></td><td>- žim - 3 <b>0,00271</b></td><td>- Zato što <b>0,00143</b></td><td>- Zato što ne - <b>0,00263</b></td></tr> <tr> <td>- m - 45 <b>0,08507</b></td><td>- im - 11 <b>0,00994</b></td><td>- 3 <b>0,00143</b></td><td>ne - 3 <b>0,00263</b></td></tr> <tr> <td>- u - 58 <b>0,10964</b></td><td>- m u - 5 <b>0,00452</b></td><td>- ěto ne - <b>0,00143</b></td><td>- što ne - <b>0,00263</b></td></tr> <tr> <td>- u - 101 <b>0,19093</b></td><td>- u - 32 <b>0,02891</b></td><td>3 <b>0,00143</b></td><td>bih - 3 <b>0,00263</b></td></tr> <tr> <td>- s - 121 <b>0,22873</b></td><td>- u s - 10 <b>0,00903</b></td><td>- nebih - <b>0,00143</b></td><td>- nebih - <b>0,00263</b></td></tr> <tr> <td>- sn - 9 <b>0,01701</b></td><td>- sin - 3 <b>0,00271</b></td><td>3 <b>0,00143</b></td><td>Izdržala - <b>0,00263</b></td></tr> <tr> <td>- no - 53 <b>0,10019</b></td><td>- sno - 8 <b>0,00723</b></td><td>- bih - <b>0,00143</b></td><td>3 <b>0,00263</b></td></tr> <tr> <td>- - - 49 <b>0,03046</b></td><td>- nov - 10 <b>0,00903</b></td><td>- izdržala - <b>0,00143</b></td><td>- - - 3 <b>0,00263</b></td></tr> </tbody> </table>		Bigrami: 39	Trigrami: 42	Bigrami riječi: 9	Trigrami riječi: 8	- Za - 61 <b>0,11531</b>	- Zat - 7 <b>0,00632</b>	- Zato i - <b>0,00143</b>	- Zato i bježim - <b>0,00263</b>	- at - 40 <b>0,07561</b>	- ato - 7 <b>0,00632</b>	- 3 <b>0,00143</b>	- i bježim - <b>0,00263</b>	- to - 45 <b>0,08507</b>	- to - 27 <b>0,02439</b>	- i bježim <b>0,00143</b>	u - 3 <b>0,00263</b>	- o - 139 <b>0,26726</b>	- o i - 10 <b>0,00903</b>	- 3 <b>0,00143</b>	- bježim u - <b>0,00263</b>	- i - 65 <b>0,12287</b>	- i - 33 <b>0,02981</b>	- bježim u <b>0,00143</b>	slove - 3 <b>0,00263</b>	- b - 62 <b>0,11720</b>	- ib - 8 <b>0,00723</b>	- 3 <b>0,00143</b>	- u slove <b>0,00263</b>	- bj - 3 <b>0,00567</b>	- bj - 3 <b>0,00271</b>	- 3 <b>0,00143</b>	Zato - 3 <b>0,00263</b>	- je - 151 <b>0,28544</b>	- bje - 3 <b>0,00271</b>	- anove <b>0,00143</b>	- anove <b>0,00263</b>	- ež - 7 <b>0,01323</b>	- jež - 5 <b>0,00452</b>	Zato - 3 <b>0,00143</b>	Zato što - <b>0,00263</b>	- ži - 17 <b>0,03214</b>	- eži - 5 <b>0,00452</b>	- Zato što <b>0,00143</b>	3 <b>0,00263</b>	- im - 21 <b>0,03970</b>	- žim - 3 <b>0,00271</b>	- Zato što <b>0,00143</b>	- Zato što ne - <b>0,00263</b>	- m - 45 <b>0,08507</b>	- im - 11 <b>0,00994</b>	- 3 <b>0,00143</b>	ne - 3 <b>0,00263</b>	- u - 58 <b>0,10964</b>	- m u - 5 <b>0,00452</b>	- ěto ne - <b>0,00143</b>	- što ne - <b>0,00263</b>	- u - 101 <b>0,19093</b>	- u - 32 <b>0,02891</b>	3 <b>0,00143</b>	bih - 3 <b>0,00263</b>	- s - 121 <b>0,22873</b>	- u s - 10 <b>0,00903</b>	- nebih - <b>0,00143</b>	- nebih - <b>0,00263</b>	- sn - 9 <b>0,01701</b>	- sin - 3 <b>0,00271</b>	3 <b>0,00143</b>	Izdržala - <b>0,00263</b>	- no - 53 <b>0,10019</b>	- sno - 8 <b>0,00723</b>	- bih - <b>0,00143</b>	3 <b>0,00263</b>	- - - 49 <b>0,03046</b>	- nov - 10 <b>0,00903</b>	- izdržala - <b>0,00143</b>	- - - 3 <b>0,00263</b>
Bigrami: 39	Trigrami: 42	Bigrami riječi: 9	Trigrami riječi: 8																																																																										
- Za - 61 <b>0,11531</b>	- Zat - 7 <b>0,00632</b>	- Zato i - <b>0,00143</b>	- Zato i bježim - <b>0,00263</b>																																																																										
- at - 40 <b>0,07561</b>	- ato - 7 <b>0,00632</b>	- 3 <b>0,00143</b>	- i bježim - <b>0,00263</b>																																																																										
- to - 45 <b>0,08507</b>	- to - 27 <b>0,02439</b>	- i bježim <b>0,00143</b>	u - 3 <b>0,00263</b>																																																																										
- o - 139 <b>0,26726</b>	- o i - 10 <b>0,00903</b>	- 3 <b>0,00143</b>	- bježim u - <b>0,00263</b>																																																																										
- i - 65 <b>0,12287</b>	- i - 33 <b>0,02981</b>	- bježim u <b>0,00143</b>	slove - 3 <b>0,00263</b>																																																																										
- b - 62 <b>0,11720</b>	- ib - 8 <b>0,00723</b>	- 3 <b>0,00143</b>	- u slove <b>0,00263</b>																																																																										
- bj - 3 <b>0,00567</b>	- bj - 3 <b>0,00271</b>	- 3 <b>0,00143</b>	Zato - 3 <b>0,00263</b>																																																																										
- je - 151 <b>0,28544</b>	- bje - 3 <b>0,00271</b>	- anove <b>0,00143</b>	- anove <b>0,00263</b>																																																																										
- ež - 7 <b>0,01323</b>	- jež - 5 <b>0,00452</b>	Zato - 3 <b>0,00143</b>	Zato što - <b>0,00263</b>																																																																										
- ži - 17 <b>0,03214</b>	- eži - 5 <b>0,00452</b>	- Zato što <b>0,00143</b>	3 <b>0,00263</b>																																																																										
- im - 21 <b>0,03970</b>	- žim - 3 <b>0,00271</b>	- Zato što <b>0,00143</b>	- Zato što ne - <b>0,00263</b>																																																																										
- m - 45 <b>0,08507</b>	- im - 11 <b>0,00994</b>	- 3 <b>0,00143</b>	ne - 3 <b>0,00263</b>																																																																										
- u - 58 <b>0,10964</b>	- m u - 5 <b>0,00452</b>	- ěto ne - <b>0,00143</b>	- što ne - <b>0,00263</b>																																																																										
- u - 101 <b>0,19093</b>	- u - 32 <b>0,02891</b>	3 <b>0,00143</b>	bih - 3 <b>0,00263</b>																																																																										
- s - 121 <b>0,22873</b>	- u s - 10 <b>0,00903</b>	- nebih - <b>0,00143</b>	- nebih - <b>0,00263</b>																																																																										
- sn - 9 <b>0,01701</b>	- sin - 3 <b>0,00271</b>	3 <b>0,00143</b>	Izdržala - <b>0,00263</b>																																																																										
- no - 53 <b>0,10019</b>	- sno - 8 <b>0,00723</b>	- bih - <b>0,00143</b>	3 <b>0,00263</b>																																																																										
- - - 49 <b>0,03046</b>	- nov - 10 <b>0,00903</b>	- izdržala - <b>0,00143</b>	- - - 3 <b>0,00263</b>																																																																										

Slika 12 – Statistika obrađenih podataka

Statistika obrađenog teksta izračunata je na način da se broj pojavljivanja pojedinih bigrama, trigrami, bigrama riječi i trigrami riječi podijeli sa zbrojem svih bigrama, trigrami, bigrama riječi i trigrami riječi koje imamo u bazi za odabrani jezik, kako je prikazano na slici 13.

$$X_{st} = \frac{nX_o}{\sum X_o}$$

Slika 13 – Matematička formula za izračun statistike

$X_{st}$  konačni je statistički izlazni podatak koji dobijemo nakon obrade;  $nX_o$  broj je pojavljivanja pojedinih bigrama, trigrami, bigrama riječi, odnosno trigrami riječi u obrađenom tekstu;  $\sum X_o$  ukupni je broj bigrama, trigrami, bigrama riječi, odnosno trigrami riječi koje smo dobili nakon obrade unešenog teksta.

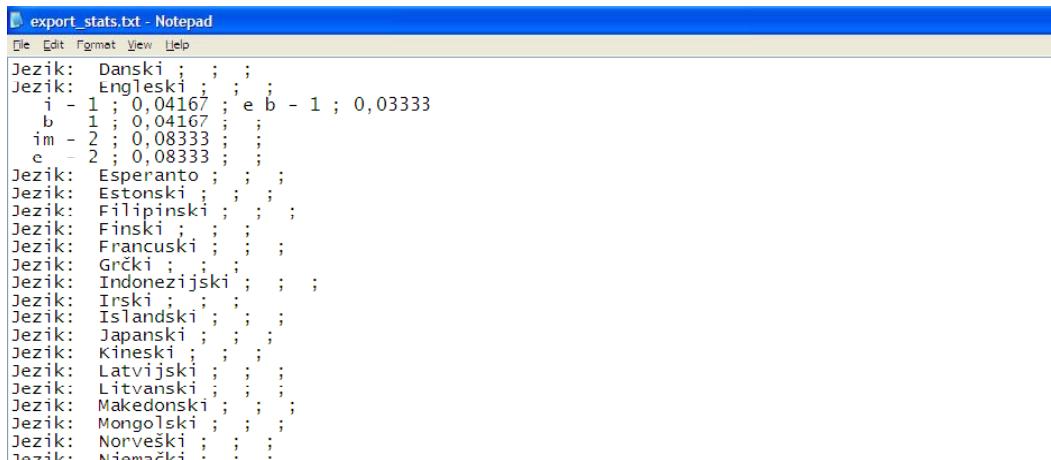
### 3.2.6. Vanjsko spremanje statističkih podataka

Osim vizualnog pregleda podataka na stranici nakon obrade, dodana je mogućnost spremanja podataka u vanjsku datoteku početnog naziva *export\_stats.txt* i to pritiskom na gumb *Spremi u .txt*. Kôd kojim spremamo statistiku prikazan je na slici 14.

```
<asp:Button ID="btnSnimi" runat="server"
    Text="Spremi u .txt"
    onclick="btnSnimi_Click" />
```

**Slika 14 – Spremanje statistike u vanjsku datoteku**

U spremljenoj datoteci nalaze se podaci iz tablice koji služe daljnjoj obradi koja se može vršiti lokalno tj. na korisnikovom računalu. Iz tog je razloga ovakvo spremanje podataka nazvano *vanjsko spremanje*. Izgled datoteke nakon spremanja na računalo i otvaranja u tekstualnom editoru prikazan je na slici 15.



```
Jezik: Danski ; ; ;
Jezik: Engleski ; ; ;
i - 1 ; 0,04167 ; e b - 1 ; 0,03333
b 1 ; 0,04167 ;
im - 2 ; 0,08333 ;
c - 2 ; 0,08333 ;
Jezik: Esperanto ; ; ;
Jezik: Estoński ; ; ;
Jezik: Filipinski ; ; ;
Jezik: Finski ; ; ;
Jezik: Francuski ; ; ;
Jezik: Grčki ; ; ;
Jezik: Indonezijski ; ; ;
Jezik: Irski ; ; ;
Jezik: Islandski ; ; ;
Jezik: Japanski ; ; ;
Jezik: Kineski ; ; ;
Jezik: Latvijski ; ; ;
Jezik: Litvanski ; ; ;
Jezik: Makedonski ; ; ;
Jezik: Mongolski ; ; ;
Jezik: Norveški ; ; ;
Jezik: Niemacki ; ; ;
```

**Slika 15 – Statistički podaci u tekstualnom obliku**

### 3.3. *Statistika teksta*

Najvažniji dio nakon obrade podataka je zapravo konačna statistika za svaki od pojedinih jezika koji se nalaze u bazi. U dalnjem tekstu objasnit ću kako i gdje doći do tih podataka, te kako dalje raspolagati dobivenim podatcima.

### 3.3.1. Statistika top30

Na stranici [www.ante.sibinj.com/Statistika.aspx](http://www.ante.sibinj.com/Statistika.aspx) moguće je pregledati globalnu statistiku obrađenih podataka koji se nalaze u bazi. Tablica na stranici prikazuje top 30 bigrama, trigramu, bigrama riječi i trigramu riječi za svaki jezik koji se nalazi u bazi poredanih na osnovu frekvencije učestalosti pojavljivanja. Pored svakog obrađenog podatka prikazan je omjer s ukupnim brojem znakova svih tekstova za pojedini jezik. Broj znakova naveden je kako bi korisnik imao referentnu vrijednost i okvir za daljnji rad, jer veća količina znakova u bazi znači veću preciznost podataka.

Izgled stranice koja prikazuje statistiku prikazan je na slici 16. Statistika ocrtava trenutno stanje baze, ovisno o broju i količini teksta koji je unešen i obrađen od početka objavljivanja aplikacije na internetu.

bigtri globalstat - top30							
Jezik:	Albanski	2225 /b./ - 2236 /zn./	372 /b.r./ - 2236 /zn./	2223 /t./ - 2236 /zn./		371 /t.r./ - 2236 /zn./	
Omjer:		2225 /b./ - 2236 /zn./	372 /b.r./ - 2236 /zn./	2223 /t./ - 2236 /zn./		371 /t.r./ - 2236 /zn./	
Statistika:	ë	0,04135	për të	0,01882	të	0,01574	transparencën e zgjedhjeve
	P	0,02517	do të	0,01613	në	0,01215	është kulmi i
	e	0,02427	nuk do	0,00806	të	0,01125	Balla ka deklaruar
	r	0,02022	e zgjedhjeve	0,00538	ër	0,00945	ka deklaruar se
	t	0,02022	ketë dhunë,	0,00538	për	0,00900	në transparencën e
	të	0,01933	të ketë	0,00538	e	0,00855	e zgjedhjeve të
	ar	0,01888	- deklaroi	0,00538	ë p	0,00765	zgjedhjeve të 28
	n	0,01708	të gjitha	0,00538	pë	0,00765	të 28 qershorit
	d	0,01618	të qenë	0,00538	në	0,00675	në seancën e
	ër	0,01483	mori vendimin	0,00538	ar	0,00630	të ketë dhunë,
	në	0,01348	i PS	0,00538	Par	0,00630	vendimin për të
	sh	0,01213	Parlamentar i	0,00538	r t	0,00585	mori vendimin për
	k	0,01213	qoftë për	0,00538	do	0,00585	Parlamentar i PS
	i	0,01124	28 qershorit	0,00538	uar	0,00585	nuk do të
	en	0,01124	dhunë, por	0,00538	Pa	0,00540	do të ketë
	a	0,01124	deklaruar se	0,00538	ka	0,00495	për të qenë
	t	0,00989	në seancën	0,00538	po	0,00495	ketë dhunë, por
	e	0,00989	seancën e	0,00538	pr	0,00495	- deklaroi deputeti

Slika 16 – Prikaz globalne statistike

Slika 17 prikazuje formulu po kojoj se izračunava frekvencija učestalosti pojavljivanja.

$$X_{st} = \frac{nX_o}{\sum X_z}$$

Slika 17 – Matematička formula za izračun frekvencijske učestalosti pojavljivanja

$X_{st}$  je konačni statistički izlazni podatak koji dobijemo nakon obrade dijeljenjem broja pojedinačnih elemenata sa ukupnim zbrojem znakova, gdje je  $nX_o$  broj

pojedinih bigrama, trigramu, bigrama riječi, odnosno trigramu riječi a  $\sum X_z$  je zbroj znakova tekstova spremljениh u bazu.

### 3.3.2. Top30 vanjsko spremanje rezultata

Statističke podatke sa stranice moguće je spremiti u vanjsku tekstualnu datoteku i to pritiskom na gumb *Snimi u .txt* koji se nalazi na istoj stranici (vidi sliku 19). Nakon pritiska na gumb dobivamo datoteku početnog naziva *export\_global.txt* u kojoj se nalaze podaci iz tablice koja je vidljiva na stranici (vidi sliku 18).

```
<asp:Button ID="btnSnimi" runat="server"
Text="Snimi u .txt"
onclick="btnSnimi_Click" />
```

Slika 18 – Spremanje statistike u vanjsku datoteku

Izgled tekstualne datoteke koja sadrži statistiku prikazan je na slici 19.

Language	Word	Frequency	Standard Deviation	Mean	Total Count
Albanian	the	118182	0,04167	1,04167	1,00000
English	the	18182	0,04167	1,04167	1,00000
Albanian	was	18182	0,04167	1,04167	1,00000
English	was	18182	0,04167	1,04167	1,00000
Albanian	of	18182	0,04167	1,04167	1,00000
English	of	18182	0,04167	1,04167	1,00000
Albanian	times	18182	0,04167	1,04167	1,00000
English	times	18182	0,04167	1,04167	1,00000
Albanian	it	18182	0,04167	1,04167	1,00000
English	it	18182	0,04167	1,04167	1,00000
Albanian	worst	18182	0,04167	1,04167	1,00000
English	worst	18182	0,04167	1,04167	1,00000
Albanian	best	18182	0,04167	1,04167	1,00000
English	best	18182	0,04167	1,04167	1,00000
Albanian	the best	18182	0,04167	1,04167	1,00000
English	the best	18182	0,04167	1,04167	1,00000
Albanian	the worst	18182	0,04167	1,04167	1,00000
English	the worst	18182	0,04167	1,04167	1,00000
Albanian	the o	18182	0,04167	1,04167	1,00000
English	the o	18182	0,04167	1,04167	1,00000
Albanian	the f	18182	0,04167	1,04167	1,00000
English	the f	18182	0,04167	1,04167	1,00000
Albanian	the ti	18182	0,04167	1,04167	1,00000
English	the ti	18182	0,04167	1,04167	1,00000
Albanian	the im	18182	0,04167	1,04167	1,00000
English	the im	18182	0,04167	1,04167	1,00000
Albanian	the me	18182	0,04167	1,04167	1,00000
English	the me	18182	0,04167	1,04167	1,00000

Slika 19 – Globalna statistika u tekstualnom obliku

Ispod gumba za spremanje u tekstualnu datoteku nalazi se i gumb *Izvezi u .xml* koji korisniku omogućava da cijelu bazu izveze kao vanjsku *xml*<sup>15</sup> datoteku. U na ovaj način izvezenoj *xml* datoteci korisnik ima cijelu bazu dostupnu na osobnom računalu. Bazu može učitati u odgovarajuće računalne aplikacije, te izvršavati željene operacije nad bazom bez potrebe ponovnog spajanja na internet. Ograničenje koje se pojavljuje kod ovakvog rada s bazom je mogućnost da se baza popuni novim podacima nakon što je korisnik spremio *xml* datoteku na računalo. Kako bi se izbjegli

<sup>15</sup> *xml* – kratica za eXtensible Markup Language - odnosno jezik za označavanje podataka

ovakvi slučajevi i kako bi korisnici u svakom trenutku imali posljednje stanje baze spremljeno na računalo, potrebno je što češće spremanje baze na računalo.

## 4. Usporedba s drugim aplikacijama

Usporedba mrežne aplikacije *bigtri* napravljena je s dvije aplikacije s različitim pristupom pojedinoj aplikaciji. Prva od ostale dvije aplikacije je mrežna aplikacija *Wolfram|Alpha* koja nudi mogućnost obrade bigrama, trigrama, bigrama riječi i trigrama riječi na sličan način. Druga aplikacija je *WordCreator* aplikacija koju korisnik pokreće lokalno s osobnog računala.

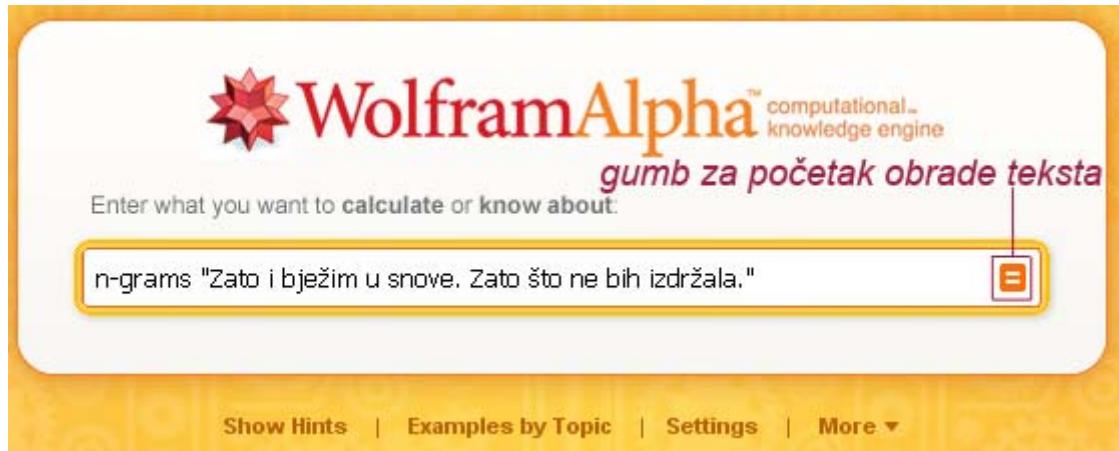
### 4.1. *Wolfram|Alpha*

*Wolfram|Alpha* je aplikacija razvijena vodeći se idejom kako bi cijelokupno sistematizirano znanje trebalo biti dostupno na što jednostavniji način. Razvoj aplikacije započeo je britanski fizičar i matematičar *Stephen Wolfram* poznat po tome što je razvio računalnu aplikaciju *Mathematica*. Zadaća aplikacije je odgovoriti korisnicima na upite koje unose u tekstualno polje, te prikazati izvore iz kojih su prikupljeni prikazani odgovori. Usporedbu aplikacije *Wolfram|Alpha*, *bigtri* i *WordCreator* započet ćemo od načina unosa teksta, te nastaviti redom do usporede dobivanja izlaznih podataka obrađenog teksta.

#### 4.1.1. Unos i interpretacija

Aplikacija *Wolfram|Alpha* koju koristimo za usporedbu dostupna je na adresi [www.wolframalpha.com](http://www.wolframalpha.com). Pristup stranici omogućen je svim internet korisnicima u bilo koje vrijeme. Tekst koji želimo obraditi unosi se na drugačiji način nego što radimo u slučaju aplikacije *bigtri*. *Wolfram|Alpha* multifunkcionalna je stranica koja odgovara na složena pitanja korisnika, te je potrebno znati točnu formu unosa podataka kako bi ga stranica pravilno interpretirala. Ispred teksta koji želimo obraditi potrebno je napisati naredbu *n-grams*, koju interpreter razumije i na taj način zna u kojem obliku želimo povratnu informaciju. Forma unosa teksta je sljedeća - ***n-grams „Zato i bježim u snove. Zato što ne bih izdržala.***“ Ovako formatiran tekst unosi se u

vodoravna komponenta za unos teksta koja se nalazi na vrhu stranice kako je prikazano na slici 20.



Slika 20 – *Wolfram/Alpha* početna stranica

#### 4.1.2. Prikaz rezultata obrade teksta

Pritiskom na gumb „=“ koji je smješten uz desni rub *tekstualne komponente*, aplikacija izvršava obradu unešenog teksta. Nakon obrade teksta, obrađeni podaci podjeljeni su i grupirani u dva zasebna dijela.

##### 4.1.2.1. Prikaz bigrama

U prvom dijelu prikazuju se samo bigrami i to bigrami znakova i bigrami riječi. Na slici 21 vidljivi su bigrami nakon obrade.

A screenshot of the WolframAlpha search results. The results section is titled "Character – level bigrams:" and displays a list of bigrams separated by vertical bars: Za | at | to | o | i | i | b | bj | je | ež | ži | im | m | u | u | s | sn | no | ov | ve | e. | . | Z | Za | at | to | o | š | št | to | o | n | ne | e | b | bi | ih | h | i | iz | zd | dr | rž | ža | a1 | la | a. Below the results is a text input field containing the sentence "n-grams \"Zato i bježim u snove. Zato što ne bih izdržala.\"". The WolframAlpha logo is visible in the bottom right corner.

Slika 21 – Bigrami dobiveni nakon obrade teksta

Slika 22 prikazuje ispis bigrama riječi.

Word – level bigrams:

Zato i | i bježim | bježim u | u snove. | snove. Zato | Zato što |  
što ne | ne bih | bih izdržala.

n–grams "Zato i bježim u snove. Zato što ne bih izdržala."

 WolframAlpha

### Slika 22 – Bigrami riječi dobiveni nakon obrade teksta

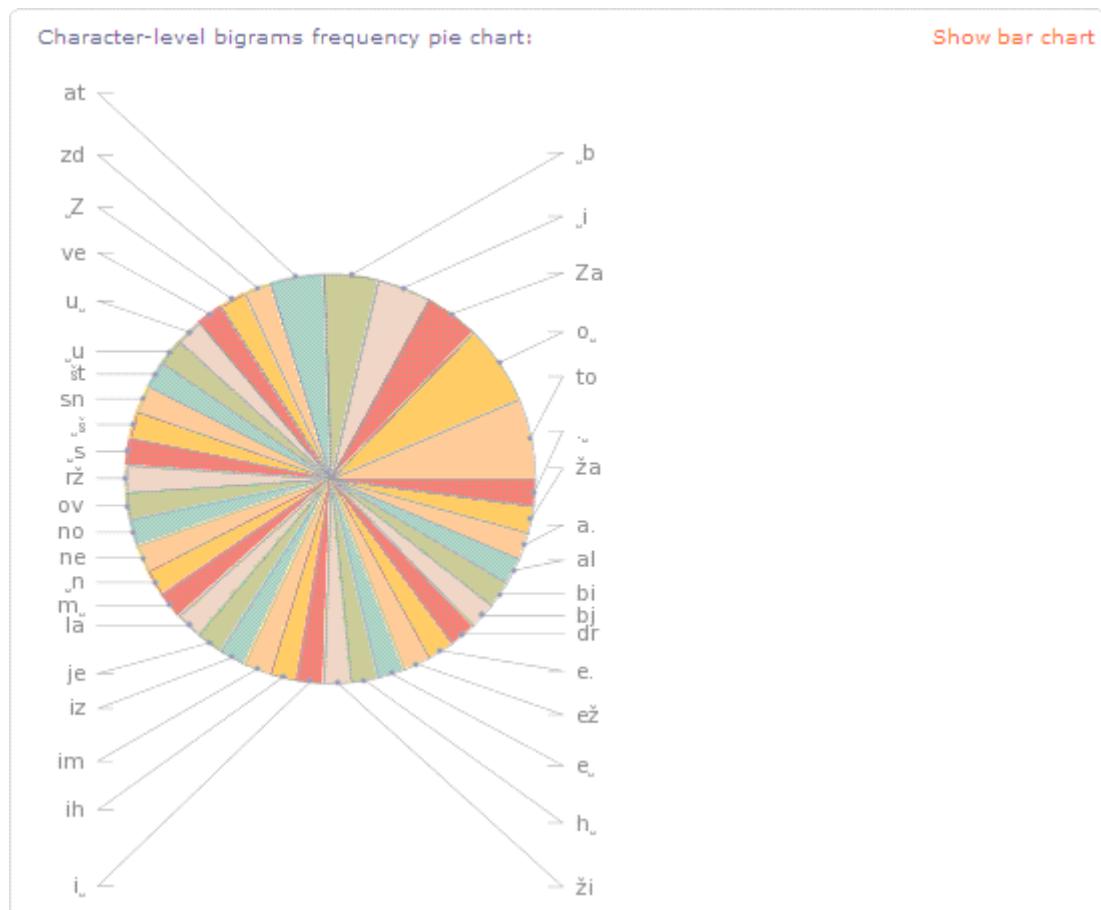
Ispod teksutalnog prikaza bigrama i bigrama riječi nalazi se detaljni grafički prikaz dobivenih podataka u obliku dijagrama. Dijagrami prikazuju broj pojavljivanja pojedinih bigrama, odnosno bigrama riječi u unešenom tekstu (njihov udio u cijelokupnom broju bigrama i bigrama riječi dobivenih iz obrađenog teksta). Gornji dijagram prikazuje broj pojavljivanja bigrama znakova (*Character-level bigrams*) u tekstu. Izgled dijagrama prikazan je na slici 23. Broj pojavljivanja bigrama riječi (*Word-level bigrams*) u tekstu prikazan je na slici 24.

Dijagrami su počentno postavljeni u *pie*<sup>16</sup> oblik, no, moguće je iste te dijagrame pregledati u *bar*<sup>17</sup> obliku pritiskom na tekst *Show bar chart*<sup>18</sup>.

<sup>16</sup> *pie* dijagram – dijagram u obliku pite

<sup>17</sup> *bar* dijagram – dijagram u obliku štapova (*štapičasti dijagram*)

<sup>18</sup> *Show bar chart* – prikaži *bar* dijagram



**Slika 23 – Dijagram bigrama obrađenog teksta**



**Slika 24 – Dijagram bigrama riječi obrađenog teksta**

#### 4.1.2.2. Prikaz trigramma

Drugi dio obrađenih podataka, koji se nalazi ispod pregleda bigrama identičan je rasporedom gornjem dijelu u kojemu su prikazani bigrami. Razlika u odnosu na gornji dio je u tome što se ovdje nalaze ispisani obrađeni trigrami i trigrami riječi. Ispis trigramma prikazan je na slici 25.

Character – level trigrams:

Zat | ato | to\_ | o\_j | \_j\_ | i\_b | \_bj | bje | jež | eži | žim |  
im\_ | m\_u | \_u\_ | u\_s | \_sn | sno | nov | ove | ve. | e\_ | \_Z |  
\_Za | Zat | ato | to\_ | o\_š | \_št | što | to\_ | o\_n | \_ne | ne\_ |  
e\_b | \_bi | bih | ih\_ | h\_i | \_iz | izd | zdr | drž | rža | žal |  
ala | la.

n–grams "Zato i bježim u snove. Zato što ne bih izdržala."

 WolframAlpha

### Slika 25 – Trigrami dobiveni nakon obrade teksta

Slika 26 prikazuje ispis trigramata riječi nakon obrade teksta.

Word – level trigrams:

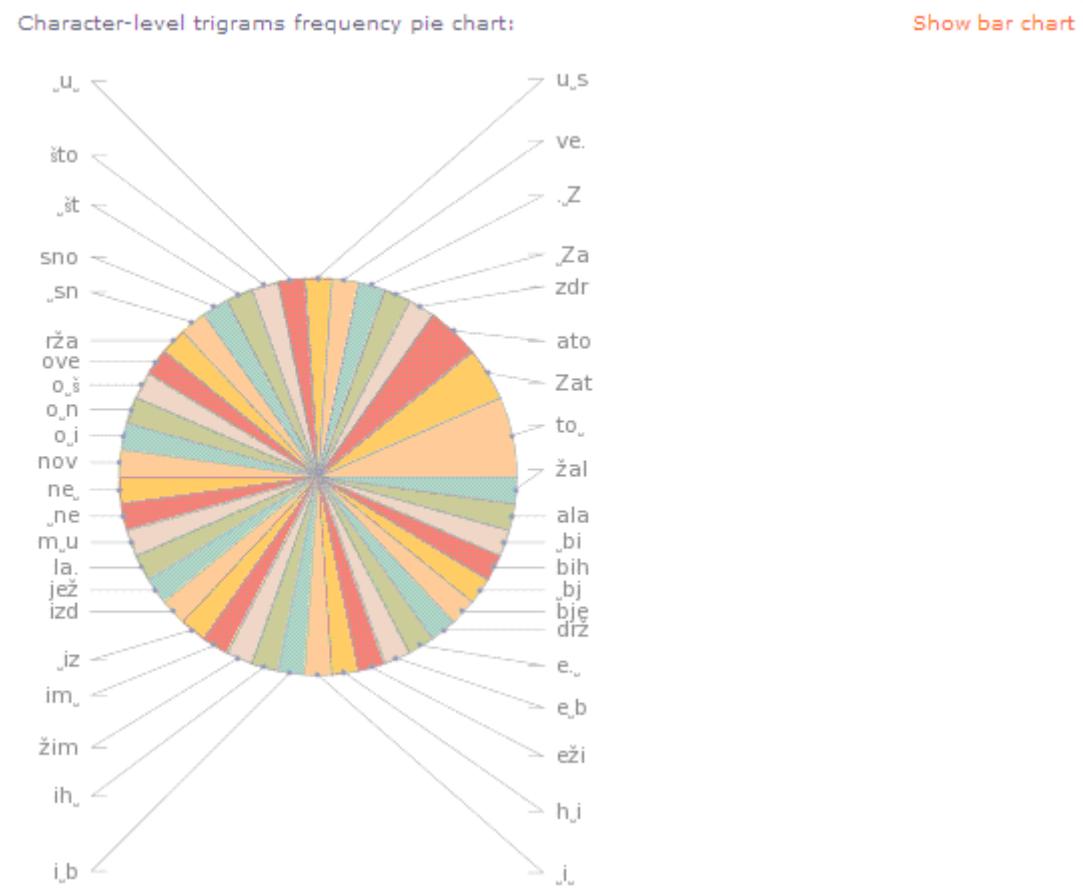
Zato i bježim | i bježim u | bježim u snove. | u snove. Zato |  
snove. Zato što | Zato što ne | što ne bih | ne bih izdržala.

n–grams "Zato i bježim u snove. Zato što ne bih izdržala."

 WolframAlpha

### Slika 26 – Trigrami riječi dobiveni nakon obrade teksta

Dijagram pojavljivanja trigramata prikazan je na slici 27, dok je na slici 28 prikazan dijagram pojavljivanja trigramata riječi. Kao što se može vidjeti, dijagrami su također početno postavljeni u *pie* oblik, te ih je moguće pregledati u *bar* obliku pritiskom na tekst *Show bar chart*



**Slika 27 – Dijagram trigrama obrađenog teksta**

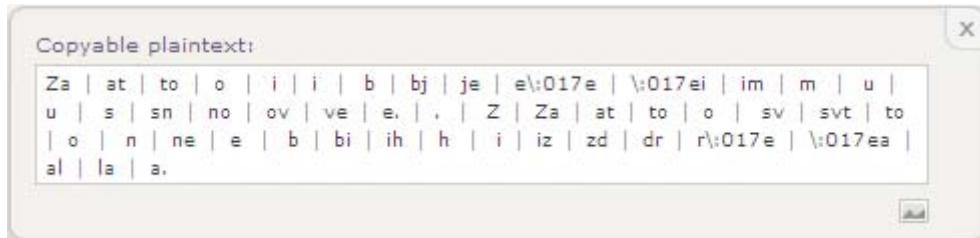


**Slika 28 – Dijagram trigrama riječi obrađenog teksta**

#### 4.1.3. Vanjska obrada rezultata obrade teksta

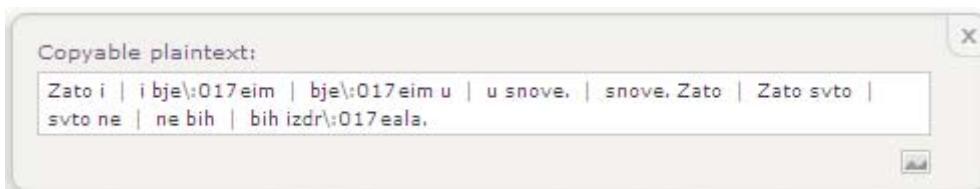
Dobivene izlazne podatke koji su prikazani na stranici nakon obrade podataka može se dobiti i u tekstualnom obliku, te kao takve koristiti lokalno za daljnju obradu. Kako bismo podatke dobili u tekstualnom obliku potrebno je kliknuti na svaku od

slika s ispisanim bigramima, odnosno trigramima. Izgled prozora s bigramima i trigramama identičan je. Razlika je jedino u znakovima koji su ispisani u njima. Tekst, odnosno znakove koje dobijemo posebno u svakom od prozora, jednostavno je označiti i kopirati te ih koristiti u dalnjem radu. Slika 29 prikazuje prozor s bigramima.



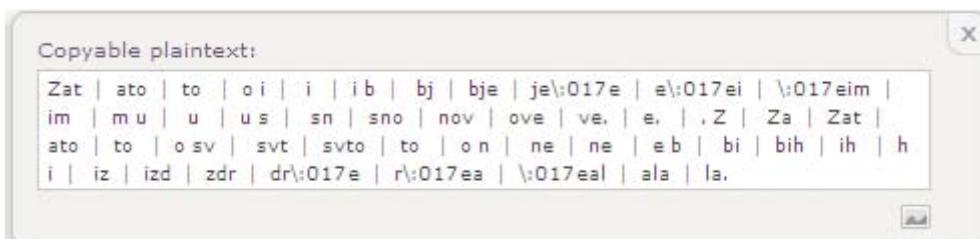
Slika 29 – Prozor s bigramima

Slika 30 prikazuje prozor koji sadrži bigrame riječi.



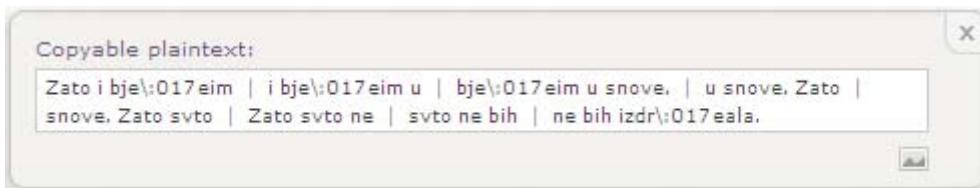
Slika 30 – Prozor s bigramima riječi

Slika 31 prikazuje prozor koji sadrži trigramme



Slika 31 – Prozor s trigramima

Slika 32 prikazuje prozor koji sadrži trigramme riječi.



Slika 32 – Prozor s trigramima riječi

Prikaz hrvatskih znakova (palatala) nije ispravan zbog korištenja UTF tablice koja podržava samo *Central European* kodiranje znakova. Korištenjem druge tablice znakova (fontova) koja u sebi ima sadržane palatale ovakva situacija bila bi izbjegnuta.

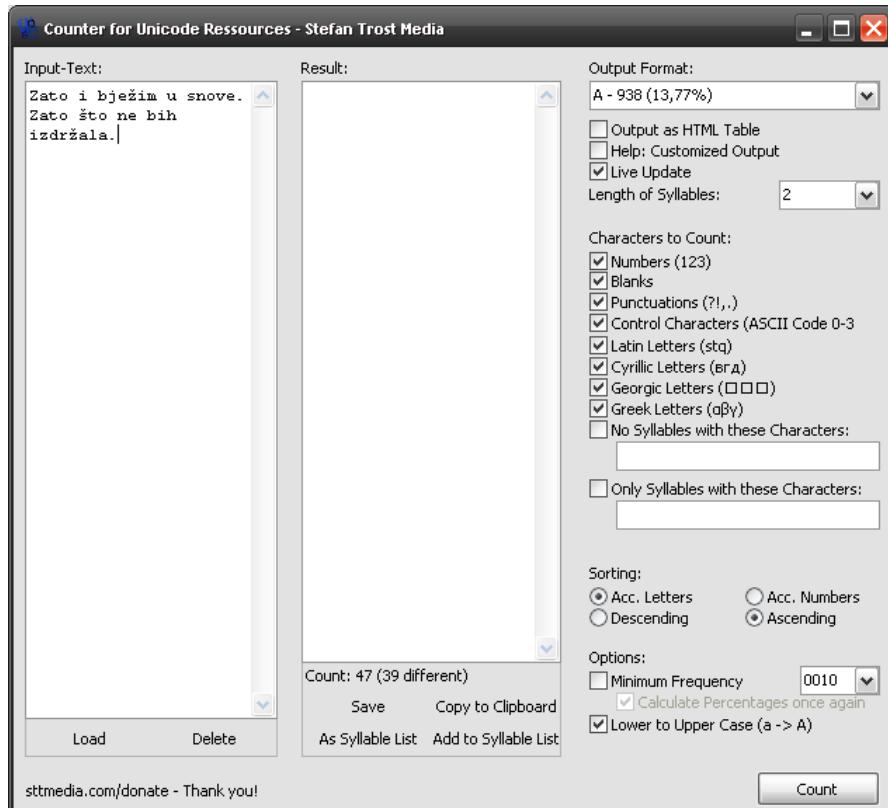
## 4.2. *WordCreator*

*WordCreator* je aplikacija koja se koristi lokalno na računalu u svrhu generiranja slučajnih nizova riječi, rečenica, te kompletnih tekstova sa subjektom i objektom. Osim generiranja, aplikacija pruža mogućnost brojanja riječi i rečenica, te prikazivanja statističkih podataka prebrojanog teksta ovisno o potrebama korisnika.

Aplikaciju *WordCreator* potrebno je snimiti na lokalni disk s internetske adrese <http://www.sttmedia.com/WordCreator>. Nakon snimanja potrebno je samo kliknuti na izvršnu datoteku i pokrenuti aplikaciju. *WordCreator* nije zahtjevna što se tiče konfiguracije računala, te ju je moguće u koristiti bez gubitka funkcionalnosti i na starijim računalima.

### 4.2.1. Unos i interpretacija

Tekst koji se obrađuje u aplikaciji *WordCreator* unosi se u lijevu komponentu za unos teksta iznad koje piše *Input-Text*. Unos teksta moguć je upisivanjem znaka po znaku putem tipkovnice, ili kopiranjem teksta s bilo koje druge lokacije. Kao primjer obrade teksta korištena je identična rečenica kao i u prethodne dvije aplikacije „*Zato i bježim u snove. Zato što ne bih izdržala.*“ Na slici 33 prikazan je pravilno unešen tekst u aplikaciju.



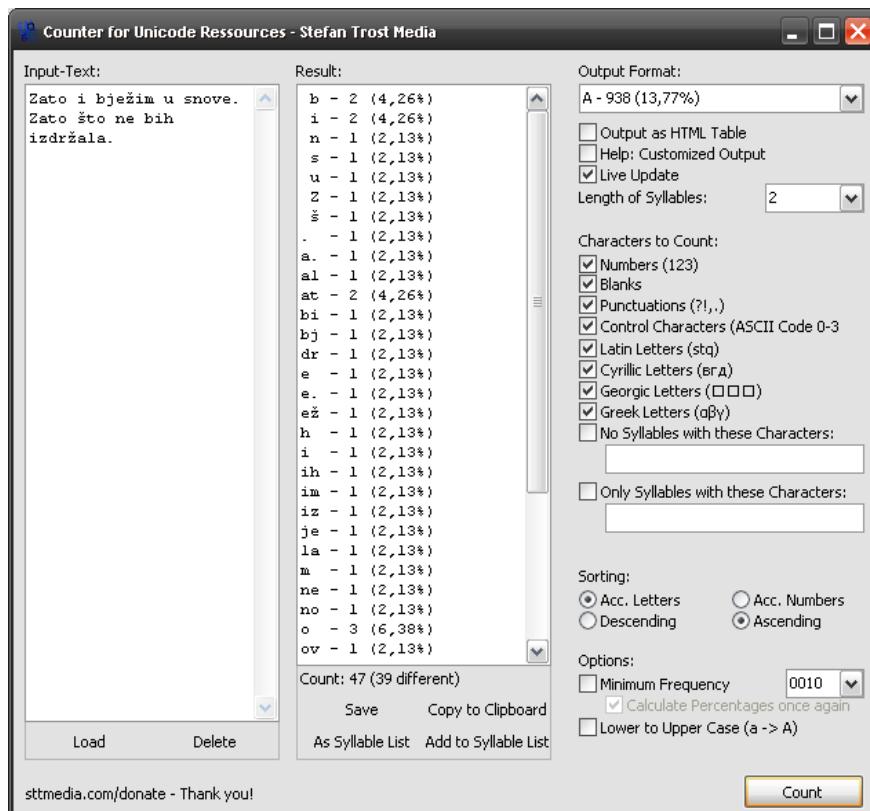
**Slika 33 – Unos teksta u *WordCreator* aplikaciju**

#### 4.2.2. Prikaz rezultata obrade teksta

Prije ispisivanja rezultata obrade teksta potrebno je izabrati na koji način želimo prikazati podatke, te koje znakove želimo brojiti, te u kojem broju (bigrame ili trigrane). Prilikom prikaza podataka moguće je izabrati želimo li rezultate prikazati samo u brojčanom obliku i to broj pojavljivanja pojedinih elemenata, ili želimo uz taj broj dodati i postotak pojavljivanja jedinice u tekstu. Duljina znakova koju prebrojavamo bira se iz padajućeg izbornika. Sve navedene opcije prikazane su uz desni rub prozora. U primjeru koji je prikazan na slici 34 iz padajućeg izbornika izabran je broj 2, što znači da želimo prikazati bigrane. Nadalje dolazimo do opcija za brojanje znakova teksta, te označavamo sve opcije prebrojavanja i to redom brojeve, prazna polja, rečenične znakove (točka, upitnik, uskličnik...), kontrolne znakove itd.

Rezultate dobivamo pritiskom na gumb *Count* koji se nalazi u donjem desnom dijelu prozora. Rezultati obrade ispisuju se u komponentu za unos teksta koja se nalazi u središnjem dijelu prozora iznad koje piše *Result*. Rezultate je moguće

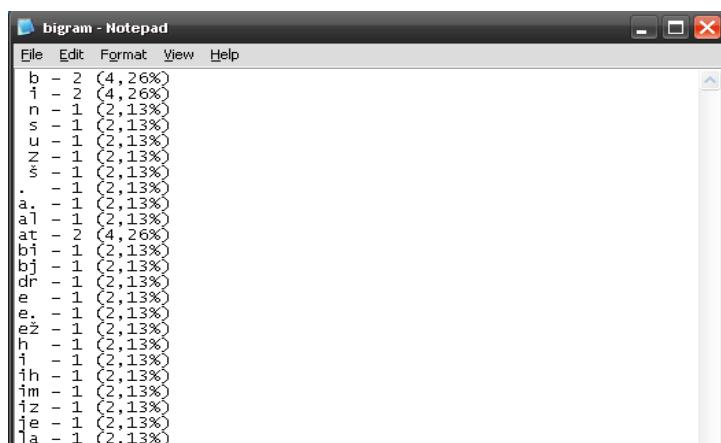
pregledavati, označavati i spremiti u memoriju računala ukoliko ih želimo koristiti u drugim aplikacijama.



Slika 34 – Rezultat obrade teksta *WordCreator* aplikacijom

#### 4.2.3. Vanjska obrada rezultata obrade teksta

Dobivene rezultate obrade teksta moguće je spremiti u vanjsku datoteku, te ih kao takve koristiti u drugim aplikacijama, ili učitati u *WordCreator* aplikaciju. Podaci se spremaju pritiskom na gumb Save. Izgled datoteke nakon spremanja na računalo i otvaranja u pripadajućoj aplikaciji prikazan je na slici 35.



Slika 35 – Rezultati obrade teksta *WordCreator* aplikacijom u tekstualnom obliku

## **5. Ocjena razlika**

### **5.1. Unos teksta**

Tekst se u sve tri aplikacije unosi preko komponente za unos teksta. Aplikacija *bigtri* ima polja za unos teksta veće nego u slučaju druge dvije aplikacije koje su spomenute u tekstu. Objasnjenje za to veoma je jednostavno.

Aplikacija *bigtri* specijalizirana je za unos veće količine teksta, te je stoga očekivano da će polje u koje se unosi tekst biti veće i preglednije.

*Wolfram|Alpha* je multifunkcionalna mrežna stranica i aplikacija koja je izrađena kako bi izvršavala mnogo različitih upita, te je stoga limit postavljen na unos teksta. Unos teksta vrši se preko komponente za unos teksta koja može primati samo jedan red znakova, što ga čini nepreglednim ako se radi o dužem tekstu koji želimo obraditi. Prilikom upisivanja teksta moramo još uzeti u obzir i naredbu koja se mora upisati prije samog teksta (*n-grams*), što dodatno može zbuniti korisnika.

*WordCreator* aplikacija ima izduženo okomito polje za unos teksta koje je također veoma pregledno prilikom unosa teksta, što ga čini preglednijim nego je to slučaj kod *Wolfram|Alpha*-e.

### **5.2. Prikaz podataka obrade teksta**

Prikaz podataka razlikuje se u svakoj od tri aplikacije. U *Wolfram|Alpha* aplikaciji tekst se prikazuje tekstualno i grafički. Grafički prikaz odnosi se na dijagrame koji su prikazani na slikama 23, 24, 27 i 28.

Aplikacija *bigtri* koristi samo tekstualni prikaz podataka kao i *WordCreator*, dok je kod ovog posljednjeg moguće izlazne podatke formatirati i u *html* kôd ukoliko je podatke potrebno objaviti na internetu.

U ovom slučaju prednost je na strani *Wolfram|Alpha*-e, jer grafički prikaz daje jasniji uvid u omjere i količinu pojedinih bigrama, odnosno trigrama u obrađenom tekstu.

### 5.3. Brzina obrade teksta

Brzinu obrade podataka možemo spomenuti u kontekstu najvažnije komponente prilikom obrade podataka. Pravilna usporedba moguća je samo ukoliko uspoređujemo aplikacije koje su dostupne na internetu, jer je nemoguće ostvariti bržu obradu nego što je to u slučaju aplikacije koja se koristi lokalno na računalu, no usporedbe radi, navest ćemo podatke i za brzinu *WordCreator* aplikacije.

U ovom segmentu aplikacija *bigtri* u velikoj je prednosti ispred *Wolfram|Alpha*-e. Mjerenja su vršena dodatkom za internetski pretraživač *Mozilla Firefox*. Prilikom pritiska gumba za obradu podataka (neovisno radi li se o *bigtri* ili *Wolfram|Alpha* aplikaciji), instalirani dodatak mjeri koliko je vremena prošlo do prikaza obrađenih podataka. Mjerenje je obavljano unutar vremenskog perioda od 60 minuta kako bi se na što bolji način izbjegla mogućnost eventualnog zagušenja internetske linije ili servera na kojemu se nalazi pojedina od aplikacija. U tablici 4 prikazana je prosječna brzina obrade podataka svih triju aplikacija. Usporedba je relevantna samo za aplikacije koje su dostupne na internetu. Jasno je vidljivo kako je brzina obrade 10.5 puta veća ukoliko koristimo *bigtri* aplikaciju.

Aplikacija	1	2	3	4	5	6	7	8	9	10	Prosjek
<i>bigtri</i>	1.208s	0.811s	0.881s	0.897s	0.936s	0.904s	0.88s	0.849s	0.817s	0.931s	0.911s
<i>WAlpha</i>	10.121s	10.31s	9.76s	9.229s	9.127s	9.995s	9.385s	9.104s	9.067s	10.274s	9.638s
<i>WCreator</i>	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s

Tablica 4 – Rezultati mjerenja brzine obrade podataka

### 5.4. Vanjska obrada rezultata

Podaci dobiveni obradom teksta koriste se ne samo kao konačna statistika nego i kao podaci kojima se želi nakon obrade koristiti lokalno na računalu.

Kako sam naveo ranije u tekstu, *Wolfram|Alpha* nema ugrađenu mogućnost statističkog pamćenja obrađenih podataka jer nije specijalizirana za takav način obrade teksta. Podatke koji su ispisani na stranici nakon obrade lako postaju dostupni korisniku, no format u kojemu ih stranica ostavlja, te nepotpunost obrade ako se radi o palatalima, uvelike otežava obradu podataka, dok aplikacija *bigtri* nema problema sa prikazivanjem palatala. Razlika je osim palatala, u prikazu leksičkih jedinica.

Kod *Wolfram|Alpha*-e se bigrami, koje sam koristio u ovom primjeru, odvajaju razmacima i znakom "|", dok aplikacija *bigtri* koristi apostrofe prije i poslije bigrama.

*WordCreator* aplikacija koristi identičan princip vanjskog prikaza podataka kao i *bigtri* aplikacija, jedino su bigrami, odnosno trigrami poslagani okomito svaki u novi red. Primjer izgleda teksta nakon obrade prikazan je u tablici 5.

Wolfram Alpha	bigtri	WordCreator
Zato i bježim u snove. Zato što ne bih izdržala.	Zato i bježim u snove. Zato što ne bih izdržala.	Zato i bježim u snove. Zato što ne bih izdržala.

Tablica 5 – Izgled teksta nakon obrade u svakoj od tri aplikacije

## 5.5. Konačna ocjena

Komparacija aplikacija donesena je na osnovu polja za unos teksta, prikaza podataka, brzine obrade, te mogućnosti vanjske obrade dobivenih podataka. Analiza je obavljena iz perspektive kako nisu svi parametri jednake vrijednosti i važnosti za korisnika.

Unos teksta, te mogućnost obrade podataka nakon obrade okarakterizirat ćemo kao najvažnije dijelove procesa obrade teksta, te će se njihov status uzimati kao presudan prilikom davanja završne ocjene. U tablici 6 prikazane su prednosti, odnosno mane koji se odnose na funkcionalnost aplikacija.

<b>Aplikacija</b>	<i>Wolfram Alpha</i>	<i>bigtri</i>	<i>WordCreator</i>
<b>Unos teksta</b>	-	+	+
<b>Prikaz podataka</b>	+	-	-
<b>Brzina obrade podataka</b>	-	+	+
<b>Vanjska obrada podataka</b>	-	+	+
<b>Konačna ocjena</b>	1 / 4	3 / 4	3 / 4

**Tablica 6 – Ocjena aplikacija**

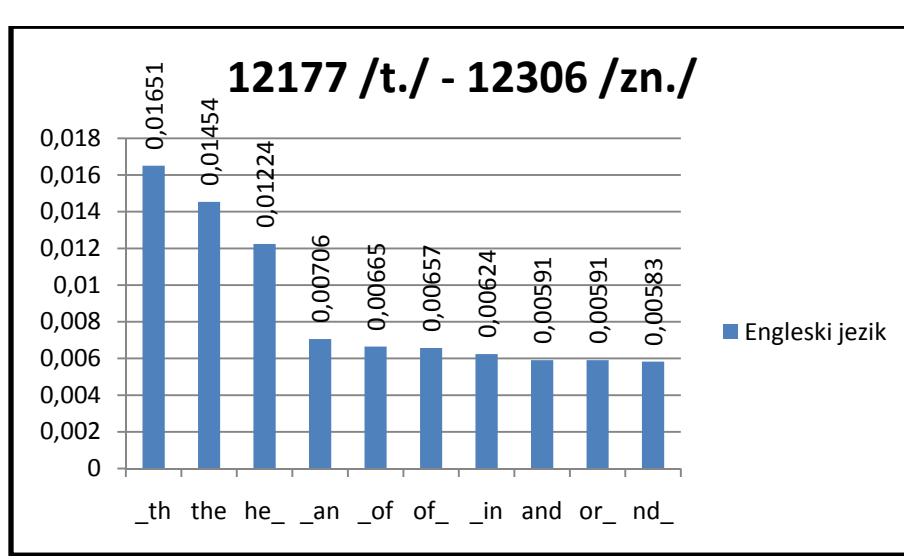
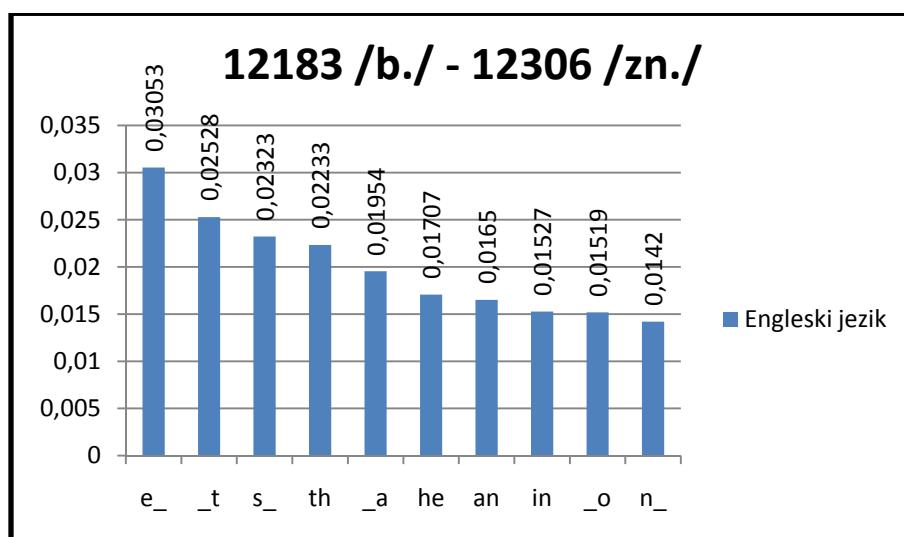
Prikaz podataka jedini je plus koji je *Wolfram|Alpha* dobila kao prednost nad *bigtri* aplikacijom. Pozitivna ocjena rezultat je prikaza grafikona, te pregledniji prikaz teksta koji je vodoravno postavljen. Vodoravni prikaz teksta prvenstveno znači kako nije potrebno dodatno pomicanje stranice prema dolje (skrolanje).

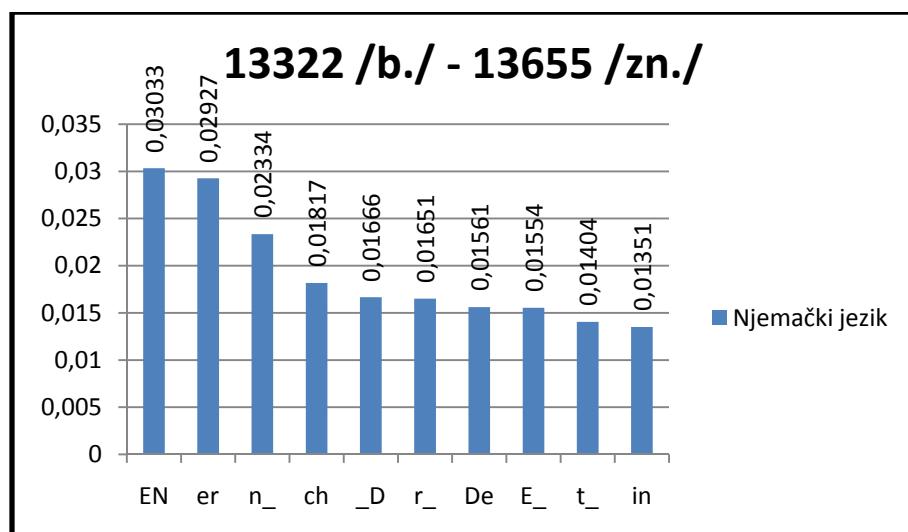
Kategorije u kojima je *bigtri* bolji od *Wolfram|Alpha*-e samorazumljive su i objašnjene ranije u tekstu. Unos teksta pregledniji je i jednostavniji za što je zaslužan veći prostor za unos teksta kao i nepostojanje dodatne naredbe kako bi aplikacija znala što s tekstrom koji je unešen. Brzina obrade podataka 10.5 puta je veća kod *bigtri* nego kod *Wolfram|Alpha*-e. Dijelom je to zbog dijagrama i ostalih grafičkih prikaza koji se mogu vidjeti na stranici. Vanjska obrada dobivenih podataka svakako je bolja i kvalitetnija nakon obrade teksta *bigtri* aplikacijom. Obrađeni tekst moguće je jednostavno spremiti na računalo u tekstualnu datoteku koju aplikacija sama generira, što je veoma jednostavno i pregledno, dok kod *Wolfram|Alpha*-e to nije slučaj.

Konačna ocjena pokazuje kako ipak specijalizirana aplikacija kao što je *bigtri* odnosi svojevrsnu pobjedu i pokazuje se kao kvalitetnije rješenje u obradi jezika i raščlanjivanje na bigrame, bigrame riječi, trigrame i trigrame riječi. *WordCreator* aplikaciju nije bilo moguće konačno pozitivno ocjeniti jer nema mogućnost obrade bigrama riječi ili trigrama riječi.

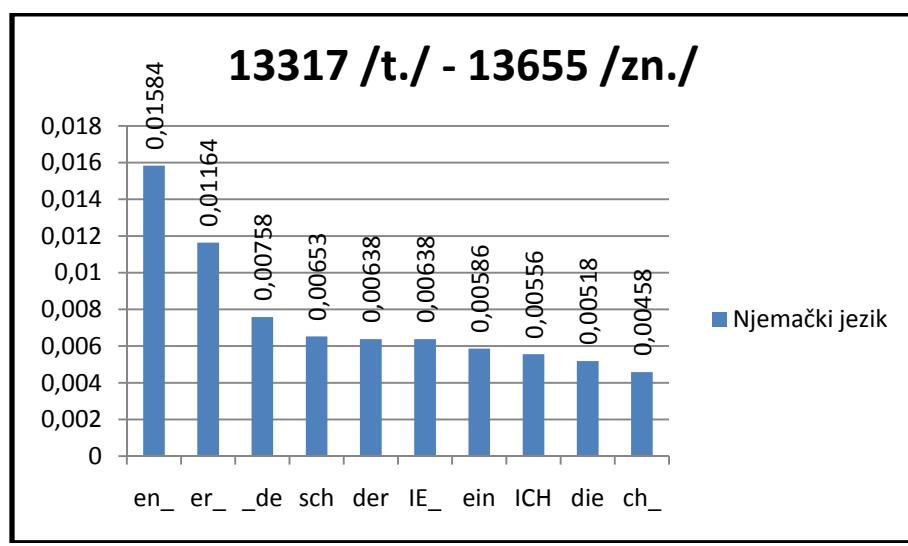
## 6. Rezultati pojavljivanja bigrama i trigrama

Kao primjer navedeni su rezultati pojavljivanja bigrama i trigrama za hrvatski, engleski i njemački jezik. Tekst koji je referentan za statistički prikaz kopiran je sa mrežnih izdanja novina za svaki jezik. U slučaju hrvatskog jezika tekst je kopiran sa internet stranice [www.jutarnji.hr](http://www.jutarnji.hr). Tekst na engleskom jeziku kopiran je sa [www.thesun.co.uk](http://www.thesun.co.uk), dok je tekst na njemačkom jeziku kopiran sa [www.bild.de](http://www.bild.de).

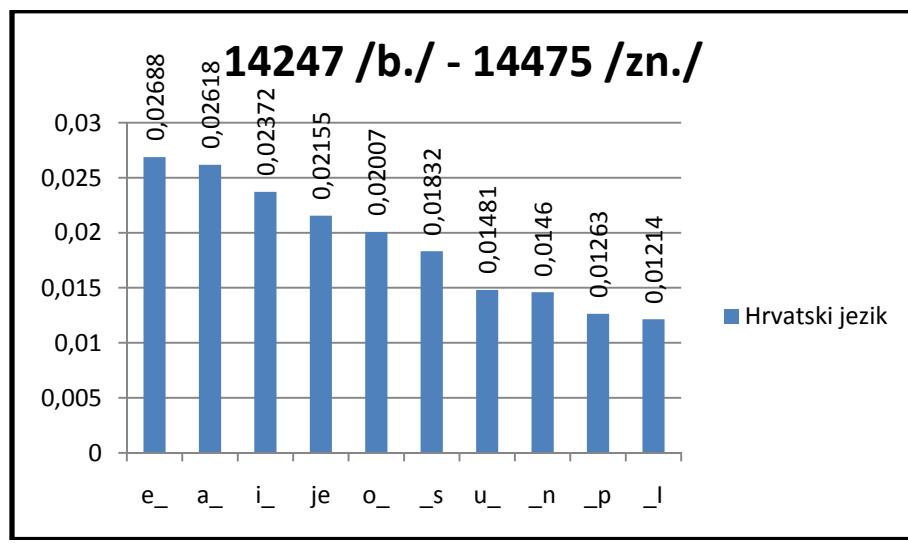




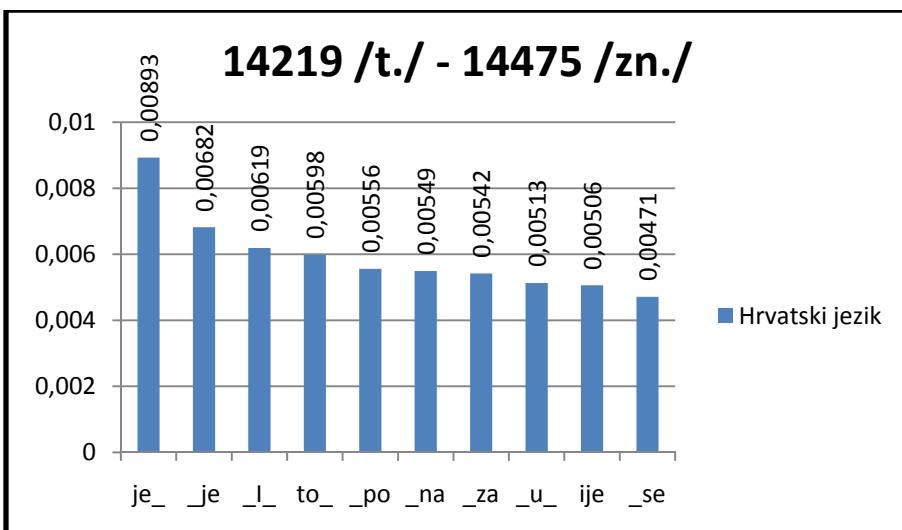
**Dijagram 3 – Njemački jezik – bigrami**



**Dijagram 4 – Njemački jezik – trigrami**



**Dijagram 5 – Hrvatski jezik – bigrami**



**Dijagram 6 – Hrvatski jezik – trigrami**

## 7. Zaključak

Izrada aplikacije koja se nalazi pred nama rezultat je, kako sam naveo u uvodu, želje za postojanjem aplikacije ovakvoga tipa na hrvatskom jeziku. Od početka izrade cilj je bio pojednostaviti korištenje i funkcionalnost aplikacije, pri čemu sam naišao na nekoliko problema i poteškoća sa spremanjem podataka u bazu i pravilnim označavanjem, te identificiranjem podataka vezanih uz jezik i tekst koji se spremaju u tablice. Uz pomoć literature i mentora, aplikacija je s vremenom dobila punu funkcionalnost, te je kasnije dodana i mogućnost izdvajanja bigrama riječi, te trigramu riječi, kao komponenti koje su prisutne u jeziku, no ne analiziraju se na adekvatan način.

Rezultati koji su dobiveni postupnom nadogradnjom programa bili su ohrabrujući, te je odlučeno kako bi bilo dobro pružiti korisnicima punu podršku i omogućiti im da mogu kompletну bazu imati na svom računalu, tako što je omogućen *izvoz* baze u *xml* formatu. Koristeći napredne razvojne alate za izradu internetskih aplikacija kreirana je aplikacija spremna udovoljiti svim zahtjevima internet standarda, kako sigurnosnih, tako i dizajnerskih.

Naposlijetku, statistički rezultati koji su dobiveni obradom bigrama i trigramu govore da je aplikacija u potpunosti referentna i da svi rezultati mogu biti korišteni u analizi jezika kakva je potrebna na najvišem znanstvenom nivou. Vjerujem da će aplikacija koju imamo pred sobom pomoći onima koji se bave ovakvom analizom jezika.

## 8. Popis literature

1. Jurafsky, Daniel & Martin, James H. : *Speech and Language Processing: An introduction to natural language processing, computational linguistics, and speech recognition*, Prentice-Hall, 2000.
2. Manning, Christopher D. & Schütze, Hinrich : *Foundations of Statistical Natural Language Processing*, MIT Press, 1999.
3. Rabiner, Lawrence R. & Juang, B. H. : *An Introduction to Hidden Markov Models*, IEEE ASSP Magazine: 4 – 16, 1986.
4. Rabiner Lawrence: *A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition*, Proceedings of the IEEE 77 (2): 257 – 286, 1989.
5. Slonneger Kenneth & Kurtz Barry L. : *Formal Syntax and Semantic of Programming Languages*, Addison – Wesley Publishing Company Inc., 1995.

Popis internet referenci:

- Wikipedia, n-gram, dostupno na Internet adresi <http://en.wikipedia.org/w/index.php?title=N-gram&oldid=367088316> , 08.06.2010.
- Wikipedia, Skriveni Markovljev model, dostupno na Internet adresi [http://en.wikipedia.org/w/index.php?title=Hidden\\_markov\\_model&oldid=175069346](http://en.wikipedia.org/w/index.php?title=Hidden_markov_model&oldid=175069346) , 08.06.2010.
- Wikipedia, Claude Shannon, dostupno na Internet adresi [http://en.wikipedia.org/w/index.php?title=Claude\\_Shannon&oldid=366047936](http://en.wikipedia.org/w/index.php?title=Claude_Shannon&oldid=366047936) , 08.06.2010.
- Jutarnji List, dostupno na Internet adresi <http://www.jutarnji.hr> , 14.05.2010.
- The Sun, dostupno na Internet adresi <http://www.thesun.co.uk> , 14.05.2010.
- Das Bild, dostupno na internet adresi <http://www.bild.de> , 14.05.2010.