

KAUNO TECHNOLOGIJOS UNIVERSITETAS
INFORMATIKOS FAKULTETAS
PROGRAMŲ SISTEMŲ INŽINERIJOS STUDIJŲ PROGRAMA

EIMANTAS ŽLABYS

AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE
SUDARYTŲ ŠABLONŲ PAGRINDU TYRIMAS

Magistro darbas

Vadovas
Prof. dr. R. Butkienė

KAUNAS, 2014

KAUNO TECHNOLOGIJOS UNIVERSITETAS
INFORMATIKOS FAKULTETAS
PROGRAMŲ SISTEMŲ INŽINERIJOS STUDIJŲ PROGRAMA

EIMANTAS ŽLABYS

AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE
SUDARYTŲ ŠABLONŲ PAGRINDU TYRIMAS

Magistro darbas

Vadovas
Prof. dr. R. Butkienė
2014-05-

Recenzentas
doc. dr. T. Blažauskas
2014-05-

Atliko:
IFM-2/2 gr. studentas
Eimantas Žlabys
2014-05-

KAUNAS, 2014

AUTENTIŠKUMO PATVIRTINIMAS

Patvirtinu, kad įteikiamas baigiamasis magistro darbas „Automatinis frazių atpažinimas tekste sudarytų šablonų pagrindu tyrimas“:

1. Autoriaus atliktas savarankiškai, jame nėra pateikta kitų autorių medžiagos kaip savos, nenurodant tikrojo šaltinio.
2. Nebuvo to paties autoriaus pristatytas ir gintas kitoje mokymo įstaigoje Lietuvoje ar užsienyje.
3. Nepateikia nuorodų į kitus darbus, jeigu jų medžiaga nėra naudota darbe.
4. Pateikia visą naudotos literatūros sąrašą.

(studento vardas, pavardė)

(parašas)

(data)

SANTRAUKA

Šiame darbe pateikiama teksto analizės bei daiktvardinių ir veiksmažodinių frazių radimo ir paieškos analizė. Analizuojami algoritmai kurių pagalba būtų galima atpažinti frazes. Algoritmai kurie sugebėtų surasti veiksmažodines bei daiktvardines frazes lietuviškame tekste. Taipogi analizuojami sprendimai bei algoritmai kurių pagalba galima būtų sudaryti šablonus, kurių dėka būtų įmanoma ieškoti frazių. Atlikus išsamią algoritmų bei esminių sprendimų analizę buvo nustatyta, jog nėra algoritmo kuri lietuviškame tekste, naudodamas vartotojo sudarytus šablonus, galėtu aptikti veiksmažodines ir daiktvardines frazes. Darbo tikslas - ištirti daiktavardinių bei veiksmažodinių frazių atpažinimą lietuviškame tekste. Siekiamas rezultatas – daiktavardinių bei veiksmažodinių frazių šablonų sudarymo programos prototipas bei algoritmo, kuris galėtu ieškoti frazių lietuviškame tekste naudojant sudarytus šablonus prototipas.

Suprojektuotas algoritmas leidžia aptikti įvairios struktūros daiktvardines ir veiksmažodines frazes. Leidžia surasti tam tikro norimo ilgio frazes. Suprojektuota frazių paieškos sistema leidžia laisvai sudaryti šablonus su norimomis detalėmis (kalbos dalis, linksnis, giminė, skaičius ir t.t.). Eksperimento metu buvo pasirinkti atsitiktiniai tekstai kuriuose algoritmas, pagal tam tikrus sistemos sudarytus šablonus, aptiko frazes. Algoritmo surastų frazių algoritmo pagrindžiamumui, buvo atliktas eksperimentas su asmenimis, kurie patys rankiniu būdu atliko frazių paiešką tuose tekstuose. Eksperimento metu buvo palygintos frazės rastos tiek eksperto tiek sistemos. Sistemos rastų frazių skaičius ne ką skiriasi nuo eksperto rastų, tačiau sistemoje sudaryti šablonai nepilnai rado visas frazes, o jei rado tai jos būdavo nepilnos. Eksperimento metu nustatyta, jog turint gerą bei tinkamą šabloną galima aptikti geresnes bei tinkamesnes frazes.

SUMMARY

This paper provides an analysis of the text, and nouns and verbal expressions of detection and analysis of the search. Analyzing algorithms which make it possible to identify phrases. The algorithms are able to find the verb and noun Lithuanian phrases in the text. We also analyzed solutions and algorithms which permit to create templates that make it possible to search for phrases. After a detailed algorithms, and key decision- analysis it was found that there is no algorithm that Lithuanian text using user -defined templates to detect noun and verb phrases. The aim - to explore the noun and verbal phrases recognition Lithuanian text. To achieve the desired result - noun and verbal phrases in the template of the program of the prototype and the algorithm is to search for phrases in Lithuanian text using templates consisting of a prototype.

Designed by the algorithm allows to detect different types of noun and verb phrases. Let's find some desired length phrases. Allows you to easily create templates with the desired parts (parts of speech, case, gender, number, etc.) . The experiment was randomly chosen texts with the algorithm in a system- defined templates found phrases. Algorithm find phrases algorithm experiment was carried out with individuals who have to manually search for phrases by those texts. During the experiment, the algorithm has been successfully implemented and applied to the analysis of the text.

TURINYS

| | |
|---|----|
| Lentelių sąrašas | 9 |
| Paveikslų sąrašas | 10 |
| Terminų ir santrumpų žodynas | 12 |
| Įvadas | 13 |
| 1. Daiktavardinių ir veiksmažodinių frazių atpažinimo tekste analizė | 14 |
| 1.1. Analizės tikslas | 14 |
| 1.2. Tyrimo objektas, sritis, problema | 14 |
| 1.2.1. SBVR standartas | 15 |
| 1.2.2. Frazių atpažinimo procesas | 17 |
| 1.2.3. Daiktavardinių ir veiksmažodinių frazių tipinės struktūros | 18 |
| 1.3. Vartotojų analizė | 22 |
| 1.4. Frazių atpažinimo įrankių analizė | 22 |
| 1.5. Frazių ir terminų atpažinimo algoritmai | 27 |
| 1.5.1. Variantų atrinkimo algoritmo analizė | 28 |
| 1.5.2. Daiktavardinių frazių atpažinimo iš sakinio struktūros algoritmas | 29 |
| 1.5.3. Atpažinimas kandidatų į žodyno elementus pagal jų formas | 29 |
| 1.5.4. Frazių ir terminų radimo algoritmų analizės išvados | 30 |
| 1.6. Siekiamo sprendimo apibrėžimas | 31 |
| 1.6.1. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama | 31 |
| 1.6.2. Dalykinės srities žodyno sudarinėtojo funkciniai reikalavimai | 32 |
| 1.6.3. Dalykinės srities žodyno sudarinėtojo sistemos nefunkciniai reikalavimai | 33 |
| 1.6.3.1. Reikalavimai sistemos išvaizdai | 33 |
| 1.6.3.2. Reikalavimai panaudojamumui | 33 |
| 1.6.3.3. Reikalavimai veikimo sąlygoms | 33 |
| 1.6.3.4. Saugumas | 34 |
| 1.7. Analizės išvados | 34 |
| 2. Frazių paieškos posistemės specifikacija ir analizė | 35 |
| 2.1. Reikalavimų specifikacija | 35 |
| 2.1.1. Šablonų sudarinėjimo posistemė | 36 |
| 2.1.1.1. Šablonų sudarymo posistemės funkcijos | 36 |
| 2.1.1.2. Šablonų sudarymo posistemės nefunkciniai reikalavimai | 37 |
| 2.1.1.3. Šablonų sudarymo posistemės sekų diagramos | 38 |
| 2.1.1.4. Šablonų sudarymo posistemės veiklos diagramos | 42 |
| 2.1.2. Teksto analizės posistemė | 43 |
| 2.1.2.1. Teksto analizės posistemės funkcijos | 43 |
| 2.1.2.2. Teksto analizės posistemės sekų diagramos | 44 |
| 2.1.2.3. Teksto analizės posistemės veiklos diagramos | 50 |

| | |
|---|----|
| 3. Automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu sprendimo algoritmas..... | 51 |
| 3.1. Daiktavardinių bei veiksmažodinių frazių paieškos veikimo aprašymas | 51 |
| 3.1.1. Šablonų sudedamosios dalys..... | 53 |
| 3.1.2. Frazių paieška naudojant šablonus..... | 54 |
| 3.1.3. Algoritmo taikymo prielaidos ir situacijos..... | 56 |
| 4. Automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu realizacijos projektas..... | 57 |
| 4.1. Sistemos architektūra | 57 |
| 4.1.1. Duomenų bazės schema | 58 |
| 4.2. Frazių paieškos posistemės detalus projektas | 59 |
| 4.3. Realizacijos modelis | 68 |
| 4.3.1. Programinių komponentų architektūra | 68 |
| 4.3.2. Diegimo modelis | 69 |
| 5. Automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu Realizacija..... | 70 |
| 5.1. Realizacijos ir veikimo aprašymas..... | 70 |
| 5.1.1. Darbo pradžia..... | 70 |
| 5.1.1.1. Instaliavimas | 70 |
| 5.1.1.2. Prisijungimas | 71 |
| 5.1.1.3. Registracija | 71 |
| 5.1.1.4. Sistemos informacija | 72 |
| 5.1.2. Sistemos langas | 72 |
| 5.1.2.1. Informaciniai ženklai..... | 73 |
| 5.1.2.2. Sistemos funkciniai laukai..... | 73 |
| 5.1.2.3. Duomenų pateikimas sistemai | 73 |
| 5.1.2.4. Informacijos pateikimas | 74 |
| 5.1.3. Frazės | 75 |
| 5.1.3.1. Atrinktos frazės | 75 |
| 5.1.3.2. Žodžių surišimas..... | 76 |
| 5.1.3.3. Frazių šablonų gamintojas..... | 77 |
| 5.2. Testavimo modelis | 80 |
| 5.2.1. Testavimo apimtis | 80 |
| 5.2.2. Testavimo strategija | 80 |
| 5.2.2.1. Vienetų testavimas..... | 80 |
| 5.2.2.2. Aukštesnio lygmens testavimas..... | 81 |
| 5.2.3. Testavimo resursai | 81 |
| 5.3. Testavimo duomenys ir rezultatai | 81 |
| 5.3.1. Vardų ir pavardžių algoritmo testavimas | 81 |
| 5.3.2. Sutrumpinimų radimo algoritmo testavimas..... | 81 |
| 5.3.3. Sakinių išskaidymo algoritmas | 82 |
| 5.3.4. Frazių radimo pagal šablonus algoritmų testavimai | 82 |

| | |
|--|-----|
| 6. Eksperimentinis Automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu sistemos tyrimas | 83 |
| 6.1. Eksperimento planas | 83 |
| 6.2. Eksperimento rezultatai..... | 83 |
| 6.2.1. Sukurtos sistemos eksperimento rezultatai | 84 |
| 6.2.2. Eilinio vartotojo atlikusio eksperimentą rezultatai | 85 |
| 6.2.3. Eksperto atlikusio eksperimentą rezultatai | 85 |
| 6.3. Sprendimo savybių analizė, kokybės kriterijų įvertinimas | 86 |
| 7. Išvados | 90 |
| 8. Literatūra..... | 91 |
| 9. Priedai | 92 |
| 9.1. 1 priedas. Sistemos rastos frazės..... | 92 |
| 9.2. 2 priedas. Eilinio vartotojo rastos frazės..... | 98 |
| 9.3. 3 priedas. Eksperto rastos frazės | 100 |

LENTELIŲ SĄRAŠAS

| | |
|---|----|
| 1 Lentelė Šablonai kurie susideda iš daiktavardžių..... | 19 |
| 2 Lentelė Šablonai kurie susideda iš būdvardžio bei daiktavardžio..... | 20 |
| 3 Lentelė Šablonai kurie susideda iš dalyvio bei aplinkui einančių daiktavardžių..... | 20 |
| 4 Lentelė Šablonų pavyzdžiai kurie susideda iš veiksmažodžio ir DF..... | 21 |
| 5 Lentelė Šablonai kurie susideda iš veiksmažodžio, prielinksnio bei DF..... | 21 |
| 6 Lentelė Šablonai kurie susideda iš veiksmažodžio, veiksmažodžio bendraties bei DF..... | 22 |
| 7 Lentelė Šablonai kurie susideda iš veiksmažodžio,rieveiksmio bei DF..... | 22 |
| 8 Lentelė Teksto analizės įrankių palyginimas..... | 27 |
| 9 Lentelė Kalbos dalių morfologinės savybės..... | 53 |
| 10 Lentelė Pagrindinės duomenų bazės schemas lentelių aprašai..... | 58 |
| 11 Lentelė FindWordsBetweenQuote metodas..... | 60 |
| 12 Lentelė FindNamesSurNameMethod metodas..... | 60 |
| 13 Lentelė FindShortWord metodas..... | 61 |
| 14 Lentelė FindSimilaritiesMethod metodas..... | 61 |
| 15 Lentelė FindLongestPhrases metodas..... | 62 |
| 16 Lentelė WritePhrasesWithoutAdverb metodas..... | 62 |
| 17 Lentelė WriteUniquePhraes metodas..... | 63 |
| 18 Lentelė FindVerbPhraseMap metodas..... | 63 |
| 19 Lentelė FrazesFoundLogic1(WordsPartsOfSentenceProperty property, PhraseStructures phraseStructureTypeses) metodas..... | 64 |
| 20 Lentelė CheckPhraseStructure(WordPartsOfSentence wordPartsOfSentence, PhraseStructureTypes phraseStructureType) metodas..... | 64 |
| 21 Lentelė PhraseStryctures metodas..... | 64 |
| 22 Lentelė CheckConnectionPD(List<Tuple<WordPartsOfSentence, PhraseStructureTypes>> phraseTypeToCheck) metodas..... | 64 |
| 23 Lentelė VerbInfinitiveCheck(PhraseStructureTypes phraseStructure, WordPartsOfSentence wordPartsOfSentence)metodas..... | 65 |
| 24 Lentelė FrazesFoundLogic2metodas..... | 65 |
| 25 Lentelė Generator metodas..... | 66 |
| 26 Lentelė GetPhrasesStructures metodas..... | 66 |
| 27 Lentelė FoundLogicNoun2 metodas..... | 67 |
| 28 Lentelė Komponentų paaiškinimų lentelė..... | 69 |
| 29 Lentelė Galimi sakinių pabaigos simboliai..... | 82 |
| 30 Lentelė Sistemos ekeperimente naudojamų šablonų sąrašas..... | 84 |
| 31 Lentelė Sistemos kiekvieno šablono rastų frazių skaičius..... | 84 |
| 32 Lentelė Ne šablonų rastų razių skaičius..... | 85 |
| 33 Lentelė Sistemos rastų frazių kiekiai kiekviename tekste..... | 85 |
| 34 Lentelė Eilinio vartotojo rastų frazių skaičius kiekviename tekste..... | 85 |
| 35 Lentelė Eksperto rastų frazių skiačius kiekviename tekste..... | 85 |
| 36 Lentelė Sistemai keltų kriterijų sąrašas bei rezultatai..... | 89 |
| 37 Lentelė Nepilnų frazių lyginant su ekspertu sąrašas..... | 89 |

PAVEIKSLŲ SĄRAŠAS

| | |
|--|----|
| 1 Pav. SBVR žodyno sududenamosios dalys | 15 |
| 2 Pav. Fakto tipo schema..... | 16 |
| 3 Pav. Frazijų atpažinimo veiklos diagrama..... | 17 |
| 4 Pav. Frazijų radimo veiklos esybės | 18 |
| 5 Pav. MultiTerm Desktop 2011 [6]..... | 23 |
| 6 Pav. SynchoTerm, teksto analizavimo įrankis | 24 |
| 7 Pav. Wordfast sąsaja, kurioje rodomi atrinkti terminai | 25 |
| 8 Pav. TexNet32 sąsaja, kurioje rodomi atrinkti terminai..... | 26 |
| 9 Pav. Daiktavardinių frazių struktūra [9]..... | 29 |
| 10 Pav. Kandidatų lentelė..... | 30 |
| 11 Pav. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama | 31 |
| 12 Pav. Dalykinės srities žodyno sudarytojo sistemos veiklos diagrama | 32 |
| 13 Pav. Frazijų paieškos posistemės panaudojimo atvejų diagrama | 35 |
| 14 Pav. Šablonų sudarinėjimo posistemės panaudojimo atvejų diagrama | 36 |
| 15 Pav. P.A. „Sudaryti šabloną“ sekų diagrama | 38 |
| 16 Pav. P.A. „Redaguoti šabloną“ sekų diagrama..... | 39 |
| 17 Pav. P.A. „Priskirti šabloną grupei“ sekų diagrama | 40 |
| 18 Pav. P.A. „Gauti kategorijas ir grupes“ sekų diagrama..... | 41 |
| 19 Pav. P.A. „Sudaryti šabloną“ veiklos diagrama | 42 |
| 20 Pav. Teksto analizės posistemės panaudojimo atvejų diagrama | 43 |
| 21 Pav. P.A. „Ieškoti frazių“ sekų diagrama | 44 |
| 22 Pav. P.A. „Ieškoti vardų bei pavardžių“ sekų diagrama..... | 45 |
| 23 Pav. P.A. „Ieškoti sutrumpinimų“ sekų diagrama | 46 |
| 24 Pav. P.A. „Ieškoti sinonimų“ sekų diagrama | 47 |
| 25 Pav. P.A. „Ieškoti frazių tarp kabučių“ sekų diagrama | 48 |
| 26 Pav. „Frazijų paieška pagal šablonus“ sekų diagrama | 49 |
| 27 Pav. „Frazijų radimo algoritmo“ veiklos diagrama | 50 |
| 28 Pav. „Frazijų radimo naudojant šablonus algoritmas“ veiklos diagrama..... | 51 |
| 29 Pav. Veiksmažodinių frazių atpažinimo veiklos diagrama | 54 |
| 30 Pav. Detali daiktavardinių veiksmažodinių frazių paieškos algoritmo veiklos diagrama..... | 55 |
| 31 Pav. Dalykinės srities žodyno sudarytojo sistemos architektūrinis sprendimas..... | 57 |
| 32 Pav. Pagrindinės duomenų bazės schema | 58 |
| 33 Pav. Frazijų radimo klasių daigramą | 59 |
| 34 Pav. Frazijų radimas naudojant šablonus klasių diagrama | 59 |
| 35 Pav. Programų komponentų architektūra | 68 |
| 36 Pav. Sistemos diegimo modelis | 69 |
| 37 Pav. Instaliavimo failai | 70 |
| 38 Pav. Instaliavimo langas | 70 |
| 39 Pav. Instaliavimo progreso langas | 71 |
| 40 Pav. Prisijungimo langas | 71 |
| 41 Pav. Registracijos langas | 72 |
| 42 Pav. Sistemos informacinis langas | 72 |
| 43 Pav. Pagrindinis programos langas..... | 72 |
| 44 Pav. Kalbos dalių langas..... | 74 |
| 45 Pav. Frazijų meniu lango pasirinkimai | 75 |
| 46 Pav. Atrinktų frazių pagrindinis langas | 75 |
| 47 Pav. Žodžių sujungėjo langas | 76 |
| 48 Pav. Frazijų sudarinėtojo langas | 77 |
| 49 Pav. Frazijų sudarinėtojo langas | 78 |
| 50 Pav. Frazijų šablonų sąrašas | 78 |
| 51 Pav. Frazijų sudarinėtojo langas redagavimo režime | 79 |
| | 10 |

| | |
|--|----|
| 52 Pav. Frazijų šablonų sąrašas | 79 |
| 53 Pav. Eksperimentavimo procesas | 83 |
| 54 Pav. Rastų frazių kiekiai | 86 |
| 55 Pav. Visuose tekstuose rastų frazių kiekiai | 86 |
| 56 Pav. Eilinio vartotojo rastų frazių palyginimas su eksperto rastomis frazėmis..... | 87 |
| 57 Pav. Sukurtos sistemos rastų frazių palyginimas su eksperto rastomis frazėmis | 87 |
| 58 Pav. Sukurtos sistemos atpažintų frazių detalesnis pjūvis..... | 88 |

TERMINŲ IR SANTRUMPŲ ŽODYNAS

SBVR - Semantics of Business Vocabulary and Business Rules

VT – veiklos taisyklės, tam tikrus veiklos aspektus apibrėžiantis arba apribojantis teiginys, kuris reikalingas norint įvertinti veiklos struktūrą arba valdyti/įtakoti veiklą.

TSZ – Teminės srities žodynas

DF – daiktavardinė frazė

VF – veiksmažodinė frazė

IVADAS

Sistemoms ir technologijoms tobulėjant – vartotojų norai didėja. Labiau norima, jog dauguma rankinio darbo galima būtų perduoti kompiuteriams. Tokie įrankiai kurie galėtų pagelbėti vartotojams greičiau ir kokybiškiau atlikti darbus, vis tobulėja ir plečiasi į rinką.

Jau nuo seno buvo norima daugiau teksto analizės įrankių kurie galėtų analizuoti ir atlikti teksto analizę. Tokie įrankiai gali labai supaprastinti ir be abejo pagreitinti teksto analizę. Tokiems įrankiams šiuo metu yra skiriamos labai didelės lėšos ir laikas. Šių įrankių paklausa visame pasaulyje yra didelė ir panaudojimo galimybės labai plačios. Įrankiai kurie sugeba analizuoti tekstus, išskirti tekstą iš nuotraukų ar perskaityti dokumentus yra labai naudingi. Kaip vienas iš pavyzdžių šiuo metu naudojamo teksto atpažinimo įrankių gali būti panaudojamas – automobilių numerių atpažinimui ar kitur. Tokios sistemos leidžia užtikrintai visą tekstą perkelti į programas kurios gali toliau analizuoti tekstą pagal vartotojo pageidavimus.

Darbo aktualumas

Norint lietuviškame tekste rasti frazes, reikia gerai išmanyti lietuvių kalbos taisykles. Teksto analizė rankiniu būdu yra sudėtingas ir daug laiko reikalaujantis procesas. Dabartiniai algoritmai bei sistemos negali iš teksto išskirti daiktavardines bei veiksmažodines frazes, naudodamiesi vartotojo sudarytais šablonais.

Tyrimo objektas - daiktavardinių bei veiksmažodinių frazių atpažinimas lietuviškame tekste.

1. Išanalizuoti:
 - 1.1. Esančius frazių paieškos algoritmus.
 - 1.2. Esančius teksto analizės įrankius.
 - 1.3. Frazių atrinkimo procesą.
2. Suprojektuoti automatinį frazių paieškos įrankį.
3. Realizuoti automatinį frazių paieškos įrankį.
4. Atlikti automatinio frazių paieškos eksperimentinį tyrimą.

Gauti rezultatai

Atlikus eksperimentą buvo nustatyta jog sistema sugeba atpažinti daiktavardines bei veiksmažodines frazes lietuviškame tekste naudojant šablonus.

Darbo struktūra

Darbą sudaro šeši skyriai. Pirmame skyriuje aprašoma tyrimo tikslas, sritis, objektas, uždaviniai, esami sprendimai ir analizės išvados. Antrame skyriuje aprašomas frazių paieškos posistemės specifikacija ir analizė. Trečiame skyriuje aprašomas automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu sprendimo algoritmas. Ketvirtame skyriuje aprašomi automatinio frazių atpažinimo tekste sudarytų šablonų pagrindu realizacijos projekto: realizacijos, architektūros ir elgsenos modelis. Penktame skyriuje aprašomas sistemos veikimas, testavimo modelis. Šeštame skyriuje aprašomas eksperimentas ir jo rezultatas. Prieduose galima rasti frazes kurias surado bei aptiko eksperimente dalyvavę asmenys bei sistema.

1. DAIKTAVARDINIŲ IR VEIKSMAŽODINIŲ FRAZIŲ ATPAŽINIMO TEKSTE ANALIZĖ

1.1. Analizės tikslas

Analizės tikslas yra išsiaiškinti daiktavardinių bei veiksmažodinių frazių sudarymo savybes. Išnagrinėti pagrindinius daiktavardinių frazių bei veiksmažodinių frazių skirtumus ir panašumus. Analizės tikslas taipogi yra suprasti tyrimo problemą, jos priežastis ir taipogi rasti tinkamą problemos sprendimą.

1.2. Tyrimo objektas, sritis, problema

Šioje dalyje apžvelgsime pagrindinius SBVR žodyno sudarymo principus bei metodus. Detaliau paanalizuosime semantinių žodyno sudarinėjimo taisykles. Taipogi šioje dalyje bus aprašoma teksto atpažinimo principai bei eiga. Taipogi bus labiau padetalizuota daiktavardinių bei veiksmažodinių frazių atpažinimo principai bei algoritmai leidžiantys lietuviškame tekste aptikti frazes.

Šio tyrimo pagrindinė sritis yra teksto automatizavimas pritaikant naujausias technologijas ir įrankius.

Problema:

Prieš pradėdant nagrinėti tyrimo objekto problemą, būtina suprasti ką reiškia sąvoka žodynas. Žodynas - yra žodžių sąrašas su jų reikšmėmis (apibrėžimais). Žodynus sudaryti yra sunku ir užima daug laiko. Sudaryti tam tikros srities specializuotus žodynus reikia tam tikrų resursų tokių kaip: specializuotos knygos, straipsniai ar kiti dokumentai. Norint šiuos dokumentus analizuoti bei sudaryti gerą specializuotą žodyną būtina gerai išmanyti tą sritį kuriai yra sudaromas žodynas. Ankščiau ir iki šių dienų sudaryti specializuotą žodyną užtrunka daug laiko bei išieikvojama daug žmogiškųjų resursų. Norint šį procesą pagreitinti, jau yra kuriamos sistemos bei algoritmai kurie gali analizuoti pateiktus tekstus ir taipogi sudaryti specializuotus žodynus. Tačiau specializuotus žodynus sudaryti Lietuvių kalbai yra sunku bei šiuolaikiniai įrankiai nepasižymi geru bei tinkamu sudarymu.

Norint gerai bei tinkamai sudaryti specializuotą žodyną būtina žinoti pagrindinius žodynų tipus. Lietuvoje bei Užsienyje vyrauja keletas pagrindinių žodynų tipų:

- Kalbos žodynai
- Dvikalbiai žodynai
- Terminų žodynai
- Semantiniai veiklos žodynai

Kalbos žodynas - yra universalus aiškinamasis, norminamasis, plačiausiam skaitytojų ratui skirtas bendrinės kalbos kodifikacijos veikalas [2].

Dvikalbiai žodynai – verstiniai dvikalbiai, trikalbiai ar daugiakalbiai žodynai. Juose ne aiškinami kurios nors vienos kalbos žodžiai, bet „išverčiami“ į kitą kalbą arba, tiksliau sakant, pateikiami vienos kalbos žodžiams kitos ar kitų kalbų atitikmenys [3].

Terminų žodynai - tai žodynai kurios nors atskiros mokslo ar technikos šakos žodžius, daugiausia terminus, fiksuojantys žodynai. Todėl jie paprastai vadinami terminologiniais (pvz., fizikos, chemijos, tekstilės, sporto, lingvistikos ir t. t. žodynai) [4].

Semantiniai veiklos žodynai - tai žodynai, kuriuos sudaro ne tik tam tikros veiklos srities žodžiai, bet ir jų junginiai, išreiškiantys tos srities sąvokas, sąvokų junginius. Šie žodynai sudaromi prisilaikant SBVR standarto ir naudojami veiklos taisyklėms aprašyti. Žodynas ir taisyklės, kurie sudaryti prisilaikant SBVR standarto, gali būti interpretuojami ir kompiuterinių programų. Tam tikrose veiklos srityse, kur yra daug taisyklių ir jos keičiasi.

Žodynų pagrindinė sudedamoji dalis yra daiktavardžiai arba daiktavardinės frazės bei veiksmažodžiai arba veiksmažodinės frazės. Šiuos junginius tekste vartotojui rankiniu būdu yra gan sudėtinga aptikti ir tai užtrunka daug laiko.

Dėl laiko taupymo sumetimų, bei greitesnio žodynų sudarymo, daiktavardinių bei veiksmažodinių frazių radimas turi būti kompiuterizuojamas. Kadangi šiuo metūnėra algoritmo ir įrankio sugebančio lietuviškame tekste, pagal vartotojo nurodytus šablonus aptikti frazių, dėl to iškyla tokie tyrimo objektas bei sritis:

Tyrimo objektas:

Daiktavardinių bei veiksmažodinių frazių atpažinimas lietuviškame tekste.

Tyrimo sritis:

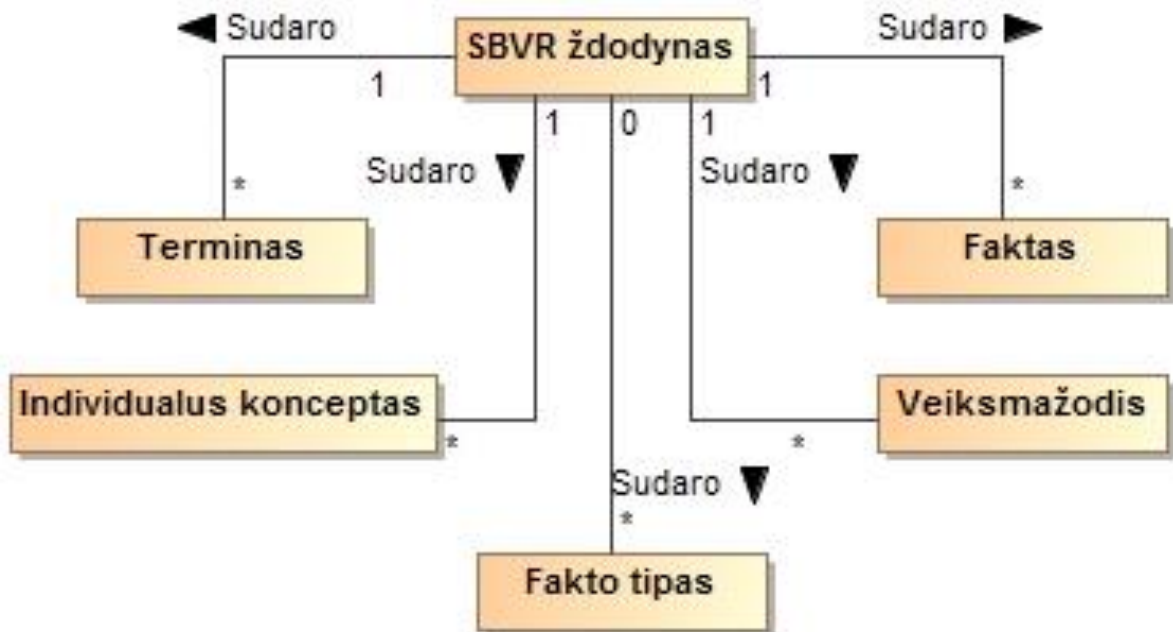
Daiktavardinių bei veiksmažodinių frazių atpažinimo lietuviškame tekste automatizavimas, automatizavimo algoritmai ir įrankiai.

1.2.1. SBVR standartas

Semantinių veiklos žodynų ir veiklos taisyklių (*Semantics of Business Vocabulary and Business Rules – SBVR*) specifikacija iš esmės apima du meta-modelius: žodyno ir taisyklių. Žodyno meta modelis tai žodynas veiklos žodynui apibrėžti(angl.), o taisyklių meta modelis - žodynas veiklos taisyklėms apibrėžti. Pastarasis naudoja ir žodyną veiklos žodynui apibrėžti.

Verslo žodynas apibrėžiamas kaip vieta, kur sukaupti visi specializuoti terminai ir koncepcijų apibrėžimai, kuriuos naudoja nagrinėjama organizacija ar bendruomenė kalbėdama ar rašydama apie savo veiklą.

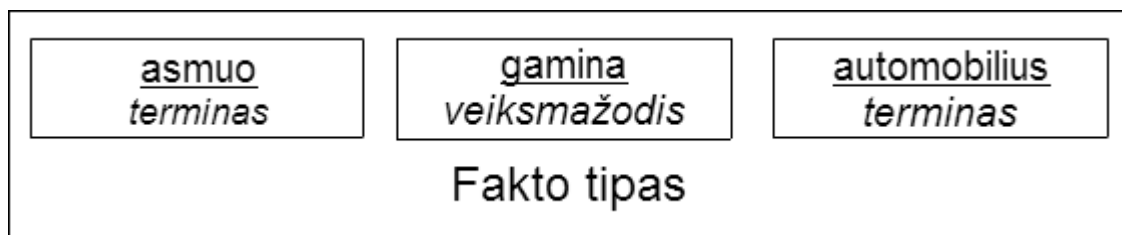
SBVR žodynus sudaro tokios pagrindinės dalys: terminai, faktai, faktų tipai ir individai (1 Pav. SBVR žodyno sududenamosios dalys). Pagrindiniai SBVR struktūrizuotos kalbos elementai kurie padetalizuoja terminų bei faktų prasmę yra terminai, individualūs konceptai, klausimai, veiksmažodžiai ir raktiniai žodžiai.



1 Pav. SBVR žodyno sududenamosios dalys

- **Terminai** – naudojami aprašyti esybėms ar rolėms, pavyzdžiui:
[universitetas](#), [studentas](#), [miestas](#), ...
- **Veiksmažodžiai** (angl. *verb*) – arba prielinksniai apibūdina terminus arba aprašo ryšius tarp terminų, pavyzdžiui:
yra, iki, sugadinta, ...
- **Individualus konceptas** – tikriniai daiktavardžiai, reiškiantys asmenis, miestus ir kitus įvardytus objektus, skaičius:
,Kaunas‘, ,Vilnius‘, ,Petras‘, ,200‘,
- **Fakto tipas** – derinys vieno fakto ir kelių terminų, pavyzdžiui:
[studentai](#) studijuoja [universitete](#), [bankas](#) suteikia [paskolą](#), ...
- **Faktai** sieja individualius konceptus arba individualius konceptus ir objektų tipus, pavyzdžiui:
[Jonas Jonaitis](#) studijuoja [KTU](#)
[Jonas Jonaitis](#) studijuoja [universitete](#)

SBVR žodyno fakto tipo sudarymą galima sudaryti pagrįdę iš dviejų elementų: termino ir veiksmažodžio (2 Pav. Fakto tipo schema). Pateiktame pavyzdyje turime sakinio iškarpa iš trijų žodžių: *asmuo gamina automobilius*.



2 Pav. Fakto tipo schema

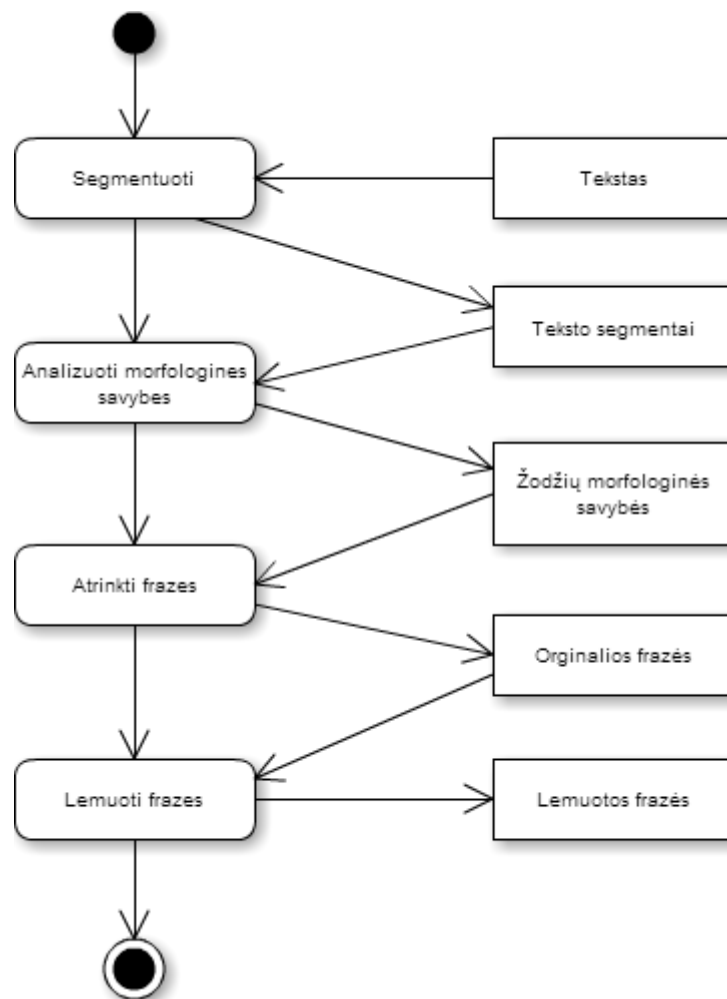
Iš šio pavyzdžio galime matyti, jog fakto tipą sudaro *veiksmažodis* kaip ryšys tarp dviejų terminų. Faktų tipo junginių gali būti įvairių, tačiau labiausiai yra naudojami:

Daiktavardis + Veiksmažodis + Daiktavardis

Dažniausiai sudarinėjant SBVR žodyną iš terminų laikomasi taisyklės, jog terminai negali būti šalia vienas kito.

1.2.2. Frazių atpažinimo procesas

Šioje dalyje apžvelgsime teksto analizavimo procesą. Plačiau panagrinėsime pagrindinius teksto analizavimo žingsnius bei atliekamus veiksmus norint rasti frazes (3 Pav. Frazių atpažinimo veiklos diagrama)

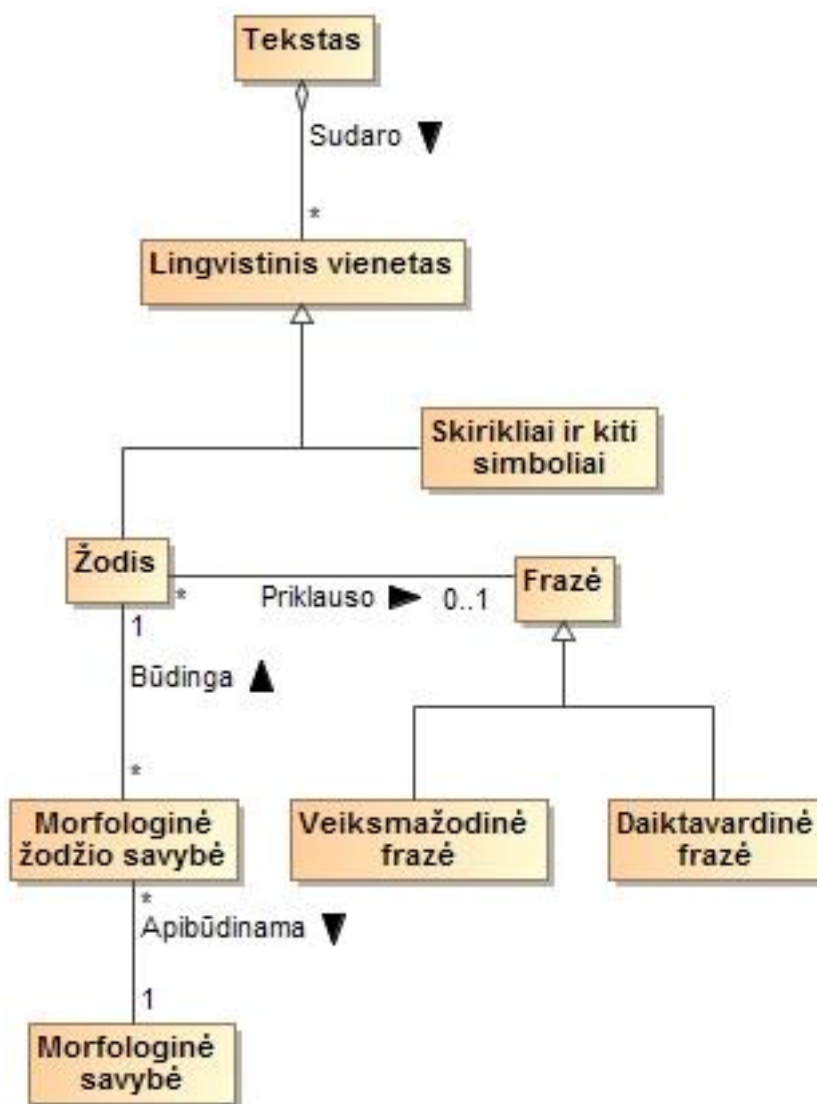


3 Pav. Frazių atpažinimo veiklos diagrama

Detaliau panagrinėsime esmines teksto analizės bei frazių radimo žingsnius kurie yra pateikti (3 Pav. Frazių atpažinimo veiklos diagrama) :

- **Segmentuoti** – šiame veiksmo atliekamas teksto segmentavimas į žodžius bei skirtukus. Šio žingsnio rezultatas yra žodžių bei lingvistinių vienetų sąrašas.
- **Analizuoti morfologines savybes** – šiame žingsnyje atliekamas kiekvieno žodžio ar skirtuko detalesnis analizavimas. Visi žodžiai bei skirtukai yra morfologiškai išanalizuojami t.y. išskiriamos *kalbos dalys*, *linksniai*, *giminė* ir t.t.
- **Frazių formavimas** – šiame žingsnyje yra ieškomos frazės, kurios gali tikti į žodyną. Šio žingsnio rezultatas yra visos rastos frazės iš teksto.
- **Frazių lemavimas** – šio žingsnio esmė yra iš rastų frazių sąrašo, pateikti visas frazes sulemuojame. Lemavimas yra žodžių ar žodžių junginių vertimas į bendrinę formą tarkime:
 - ♦ **Daiktavardžių bei būdvardžių lema** – paprastai šių kalbos dalių lema laikoma vienaskaitos vardininko forma, pvz:
gražų Kauno Technologijos Universitetą → *gražus Kauno Technologijos universitetas*
 - ♦ **Veiksmažodžių lema** – yra tiesiog bendratis, pvz:
skubėjo į → *skubėti į*

Pažvelgus į šį frazių atpažinimo procesą iš kitos perspektyvos t.y. veiklos pusės, aiškėja kitos pagrindinės teksto analizės bei žodynų esybės.



4 Pav. Frazių radimo veiklos esybės

Kaip matyti iš (4 Pav. Frazių radimo veiklos esybės) paaveiklėlo jog iš teksto pirmiausia būtina išskirti „lingvistinius vienetus“, kurie yra susideda iš dviejų dalių, t.y. žodžio bei skyriklio ir kitų simbolių. Kiekvienas žodis Lietuvių kalboje turi savo morfologines dalis, t.y. kalbos dalį, giminę, linksnį, skaičių ir t.t.. Frazę gali sudaryti daug žodžių, kadangi frazė yra žodžių junginys. Tačiau ne kiekvienas žodis gali patekti į frazės sudėtį ir žodis frazėje gali būti daugiausiai vieną kartą. Frazė – baigtinę mintį reiškiantis ir antinaciškai baigtas žodžių junginys ar posakis. Tokius junginius galima skirstyti į dvi dalis t.y. Daiktavardines frazes bei veiksmazodines frazes.

1.2.3. Daiktavardinių ir veiksmazodinių frazių tipinės struktūros

DAIKTAVARDINĖ FRAZĖ (angl. *noun phrase*) tradiciškai apibrėžiama kaip frazė, kurios pagrindinis dėmuo yra daiktavardis arba daiktavardį pakeičiantis įvardis. Iš pirmo žvilgsnio panašiai ir lietuvių kalbotyros tradicijoje apibrėžiama sąvoka „daiktavardinis junginys“, plg.: „daiktavardiniais junginiais laikomi žodžių junginiai, kurių pagrindinis dėmuo yra daiktavardis“ [5]. Tačiau kadangi tradicinėje lietuvių kalbotyroje, skirtingai nei Vakarų kalbotyroje, daiktavardiniais junginiais nelaikomi daiktavardiškieji įvardžiai ir jų junginiai (pvz., *jis, tas su geltonomis gatbanomis*), taip pat daiktavardiškai vartojami būdvardžiai (pvz., *pastarasis, aktyvioji*) bei neišplėsti daiktavardžiai (pvz., *mergaitė* ir *gėlės sakinyje mergaitė laistė gėles*), šiame darbe pasirinkta vartoti tarptautinį terminą

DAIKTAVARDINĖ FRAZĖ¹, kurio vadinami argumento poziciją užimantys daiktavardžiai, daiktavardiniai junginiai, daiktavardiškieji įvardžiai ir jų junginiai, daiktavardiškai vartojami būdvardžiai ir kt.

Pagrindinės daiktavardinių formos sudaro daiktavardis kartu einantis su būdvardžiu. Šių junginių kombinacijų sudaryti galima be galo daug, kadangi tiek daiktavardžiai tiek būdvardžiai gali kartotis daug kartų frazėje.

Frazių šablonų sudarymui yra naudojami raktiniai simboliai, kurie apibendriną šabloną:

D - daiktavardis


B – būdvardis

Dal - dalyvis

[] – privalomumas (elementas gali būti privalomas arba ne)

* - pasikartojamumas (elementas gali kartotis *n* kartų)

- Daiktavardinių frazių šablonas kurį sudaro daiktavardžiai (1 Lentelė Šablonai kurie susideda iš daiktavardžių):

$$[D]^* + D$$


1 Lentelė Šablonai kurie susideda iš daiktavardžių

| Daiktavardinės frazės struktūra | Pavyzdys |
|---------------------------------|-----------------------------------|
| D | Kaunas |
| DD | technologijos universitetas |
| DDD | Kauno technologijos universitetas |

Šios struktūros šablonas leidžia tekste aptikti bet kokio ilgio daiktavardines frazes kurios yra sudarytos iš daiktavardžių. Taipogi šis šablonas gali surasti ir vieno ilgio daiktavardinę frazę t.y. „Kaunas“ kadangi šiame šablone, pirmasis daiktavardis neprivalomas. Tačiau jei daiktavardis randamas tai jis gali kartotis *n* kartų.

¹ Iš lietuvių kalbininkų daiktavardiškai vartojamus įvardžius „daiktavardžio frazėmis“ laiko Rosinas (Rosinas 2009: 54). Mikulskas vartoja „daiktavardinės frazės“ sąvoką, kuri visiškai atitinka anglų noun phrase (Mikulskas 2006a).

- Daiktavardinių frazių šablonas kurį sudaro būdvardis bei daiktavardis. Būdvardis gali kartotis n kartų, tačiau jo gali ir nebūti daiktavardinėje frazėje (2 Lentelė Šablonai kurie susideda iš būdvardžio bei daiktavardžio):

$[B]^* + D$



2 Lentelė Šablonai kurie susideda iš būdvardžio bei daiktavardžio

| Daiktavardinės frazės struktūra | Pavyzdys |
|---------------------------------|-----------------------------|
| D | Kaunas |
| BD | naujasis kapitonas |
| BBD | Naujasis, paprastasis namas |

Šios struktūros šablonas leidžia tekste aptikti bet kokio kiekio būdvardžio ir vieno daiktavardžio frazes. Taipogi šis šablonas gali aptikti tik vieną daiktavardį kadangi būdvardis yra neprivalomas šios struktūros elementas. Žemiau pateiktame pavyzdyje matosi kaip iš sakinio galima išskirti tokio tipo frazes:

$[B]^* + D$
 ↑ ↑
 Apypigei skarelei trūko pinigų

- Daiktavardinių frazių šablonas kurį sudaro du daiktavardžiai bei per vidurį esantis dalyvis (3 Lentelė Šablonai kurie susideda iš dalyvio bei aplinkui einančių daiktavardžių):

$D + [Dal]^* + D$



3 Lentelė Šablonai kurie susideda iš dalyvio bei aplinkui einančių daiktavardžių

| Daiktavardinės frazės struktūra | Pavyzdys |
|---------------------------------|---|
| DD | Technologijos Universitetas |
| DDalD | mama mėgdžiojanti katę |
| DDalDalD | mama valganti ir besigardžiuojanti pietumis |

Šios struktūros šablonas leidžia tekste aptikti frazes kurios yra sudarytos iš dviejų daiktavardžių bei per vidurį einančio dalyvis. Detaliau panagrinėkime paskutinį lentelės (3 Lentelė Šablonai kurie susideda iš dalyvio bei aplinkui einančių daiktavardžių) pavyzdį:

$D + [Dal]^* + D$
 ↙ ↘
 mama valganti ir besigardžiuojanti pietumis

Kaip matyti iš pateikto pavyzdžio, jog tarp dviejų būdvardžių yra jungtukas „ir“, kuris turētu neleisti rasti šios frazės. Tačiau pagal nutylėjimą lietuvių kalboje yra jog prieš tai ėjęs elementas ir po jo sekantis einantis elementas yra tokios pat kalbos dalies ir tarp jų yra jungtukas „ir“, tai galima šį jungtuką ignoruoti, tačiau įtraukti į galutinę frazės vaizdą.

VEIKSMAŽODINĖ FRAZĖ (verb phrase) susideda pagrinde iš veiksmažodžio ar veiksmažodžių junginių. Viena iš pagrindinių veiksmažodinių frazių taisyklių, jog veiksmažodis privalo eiti kartu su daiktavardine fraze, tačiau daiktavardinė frazė nepriklauso veiksmažodinei frazei.

Frazių šablonų sudarymui yra naudojami raktiniai simboliai, kurie apibendriną šabloną:

V – veiksmažodis

Prl – prielinksnis

Prv –rieveiksmis

DF – daiktavardinė frazė

[] – privalomumas (elementas gali būti privalomas arba ne)

* - pasikartojamumas (elementas gali kartotis *n* kartų)

- Veiksmažodinių frazių šablonas kurį sudaro vienas veiksmažodis:

V + DF



4 Lentelė Šablonų pavyzdžiai kurie susideda iš veiksmažodžio ir DF

| Veiksmažodinės frazės struktūra | Pavyzdys |
|---------------------------------|--|
| V + DF | turi automobilį, neša maišą, bėgo namo |

- Veiksmažodinių frazių šablonas kurį sudaro veiksmažodis ir prielinksnis:


V + [Prl] + DF



5 Lentelė Šablonai kurie susideda iš veiksmažodžio, prielinksnio bei DF

| Veiksmažodinės frazės struktūra | Pavyzdys |
|---------------------------------|--|
| V + DF | turi automobilį, neša maišą, bėgo namo |
| V + Prl + DF | neša ant stogo, bėgo į urvą |


- Veiksmažodinių frazių šablonas kurį sudaro du veiksmažodžiai, bet antras veiksmažodis turi būti bendratis:

$$V + V_{(\text{bendratis})} + DF$$


6 Lentelė Šablonai kurie susideda iš veiksmažodžio, veiksmažodžio bendraties bei DF

| Veiksmažodinės frazės struktūra | Pavyzdys |
|-----------------------------------|--------------------------------------|
| $V + V_{(\text{bendratis})} + DF$ | turi būti namas, gali pamatyti krūmą |

- Veiksmažodinių frazių šablonas kurį sudarorieveiksmis ir veiksmažodis:

$$Prv + V + DF$$


7 Lentelė Šablonai kurie susideda iš veiksmažodžio,rieveiksmio bei DF

| Veiksmažodinės frazės struktūra | Pavyzdys |
|---------------------------------|--|
| $Prv + V + DF$ | smarkiai lijo kieme, sunkiai kalba meška |

1.3. Vartotojų analizė

Vartotojai atliekantys specializuotų žodynų sudarymą yra tos srities specialistai. Šios grupės vartotojų pagrindinis noras, greitai bei kokybiškai atlikti teksto analizę. Šio tipo vartotojai puikiai išmano sritį kuriai žodynas yra kuriamas bei, kur jis bus panaudojamas.

Pagrindinė vartotojų problema yra lėtas ir gan sudėtingas žodyno sudarymas bei skirtingas frazių ar faktų išskyrimas tekste. Specializuotų Lietuviškų žodynų sudarymas kompiuterių pagalba dar nėra pakankamai išvystytas, jog leistu kokybiškai išskirti daiktavardines bei veiksmažodines frazes. Vartotojai norėdami išanalizuoti didelius kiekius tekstų, privalo kaupti daug informacijos geresniems specializuotiems žodynams sudaryti.

1.4. Frazių atpažinimo įrankių analizė

Pasaulyje jau yra sukurta automatinių įrankių, kurie padeda automatiškai atrinkti terminus ir frazes. Iš jų galima sudaryti žodynus. Programos gali veikti vienoje kalboje arba gali veikti keliose kalbose. Šiame skyriuje bus apžvelgta populiariausios mokamos ir nemokamos programos skirtos ieškoti tekstuose frazių ir terminų.

Frazių paieškos įrankių analizė:

- **SDL MultiTerm Desktop S2011**
- **SynchroTerm**
- **WordFast**
- **TexNet32**

Žemiau pateikiama, detalesnė šių įrankių analizė:

SDL MultiTerm Desktop S2011

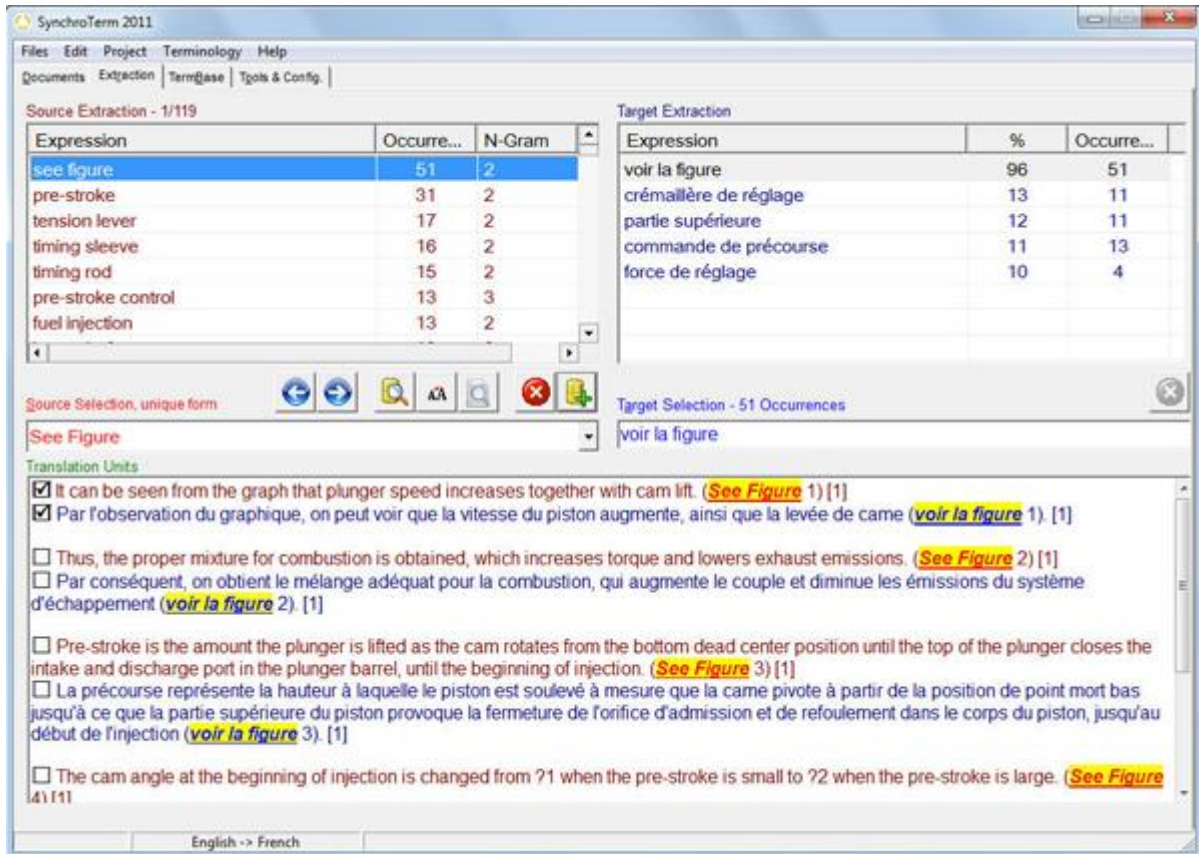
Šią mokamą (606 doleriai metams) programą sukūrė SDL kompanija ir buvo išleista į rinką 2011. Ji palaiko visas didžiąsias vakarų valstybių kalbas, norint su šia programa sėkmingai dirbti kitose kalbose reikia įdiegti papildomų priedų. SDL MultiTerm išrenka potencialius terminus, jei reikia parašo jų vertimą ir parodo galimų kandidatų sąrašą (1.pav). Sąrašą galima koreguoti ir keisti. Sutvarkytą sąrašą galima išsaugoti įvairiais formatais (XML, TXT, DOC), tai pat galima saugoti ir duomenų bazėje. Saugant į duomenų bazę galima pasirinkti ar norime saugoti į egzistuojančią ar norime sukurti naują. Galima uždėti įvairius filtrus paieškai (terminai bus ne mažiau dviejų žodžių). Programa suteikia galimybę ieškoti iš skirtingomis kalbomis parašytų tekstų [6].



5 Pav. MultiTerm Desktop 2011 [6]

SynchroTerm

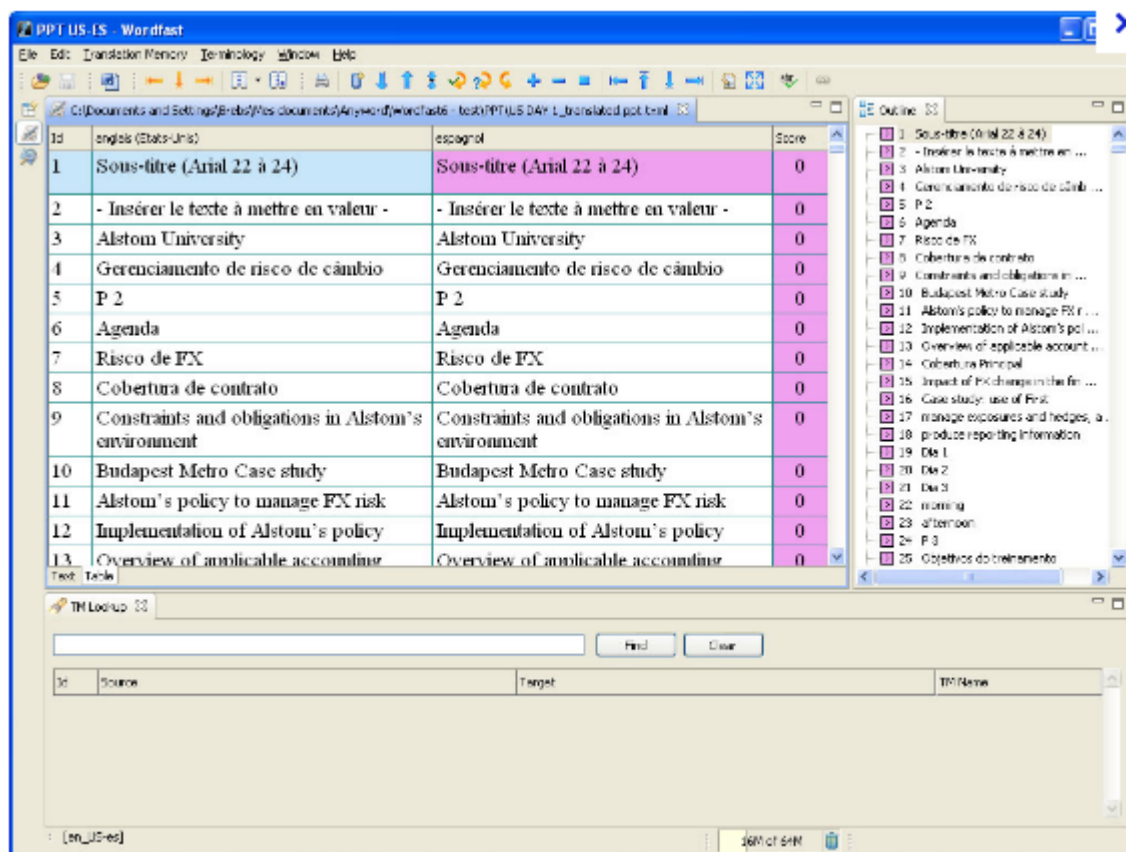
Šią mokamą (420 dolerių metams) programą išleido Terminotix kompanija. Ji palaiko didžiąsias vakarų Europos kalbas (anglų, prancūzų, ispanų ir t. t.). Galima rinkti terminus iš vieno arba kelių tekstų. Išrinktus terminus parodo sąrašė, kurį galima koreguoti (2.pav). Terminus saugoma duomenų bazėse. Duomenų bazėse esančius terminus galima koreguoti, perkelti į kitą duomenų bazę arba pašalinti iš jos. Synchterm naudoja statistinius, morfologinius ir sintaksinius algoritmus rasti terminams [7].



6 Pav. SynchoTerm, teksto analizavimo įrankis

WordFast

Šią mokamą (500 dolerių metams) programą sukūrė Wordfast LLC kompanija. Ši programa naudoja Microsoft Word kaip teksto modifikatorių. Ji veikia dauguma platformų (Windows, Mac OS, Linux). Ji palaiko visas kalbas, kurias palaiko Microsoft Word. Į jas įeina tokios kaip kinų, japonų, korėjiečių kalbas, tai pat kalbos, kuriose rašoma iš dešinės į kairę (arabų, hebrajų). Viename žodyne gali būti iki 250 tūkstančių įrašų. Tai pat į Wordfast gali apdoroti ne tik paprastu Word dokumentus, bet tai pat ir Excel'io ar Power Point. Wordfast palaiko ir HTML dokumentus. Šioje programoje galima pasirinkti kaip norima, kad atrodytų žodynas (su aprašu ar ne). Atrinktus terminus rodo Word dokumente, kuriame galima juos koreguoti. Sudarytais žodynais gali naudotis daug vartotojų per tinklą [8].



7 Pav. Wordfast sąsaja, kurioje rodomi atrinkti terminai

TexNet32

TexNet32 yra nemokama programa, kuri padeda išrinkti žodžius, frazes. Veikia tik Windows aplinkoje. TexNet32 gali nuskaityti .txt arba .rtf formatus. Tai pat failus sukurtus su TexNetF(.doc- čia ne Microsoft Word formatas ir .net). Galima apdoroti tai pat ir HTML dokumentus. Išsaugoti galima trimis formatais .rtf, .txt arba .htm. Atrinktus žodžius galima redaguoti arba ištrinti. Po koregavimo žodžius ir terminus pateikia atskirame lange. Vienu metu apdoroja tik viena dokumentą. Atrinkimas vyksta pagal kiek daugiausia simboliu gali būti frazėje (leidžiama, kad daugiausia gali būti 10 simbolių). Tai pat žiūrima kiek kartų tekste jie pasikartojo. Ieškoma tai pat ir stop žodžių (anglų kalboje tai yra *an, the*, lietuvių kalboje būtų atitikimas - *prie*).



8 Pav. TexNet32 sąsaja, kurioje rodomi atrinkti terminai.

Išanalizuotų sistemų palyginimas

Šiame skyriuje palyginsime jau aukščiau išanalizuotas sistemas. Sistemos bus įvertintos pagal šiuos kriterijus:

1. Galimybė analizuoti Lietuvišką tekstą.
2. Galimybė pateikti lemuotas frazes.
3. Galimybė iš teksto išrinkti tam tikro ilgio ar struktūros frazes.
4. Galimybė vartotojui sudaryti paieškos tekste šablonus
5. Sistema gali analizuoti: .doc, .docx, .txt, .pdf failų formatus

Teksto analizės įrankių palyginimas pateikiamas (8 Lentelė Teksto analizės įrankių palyginimas) lentelėje.

8 Lentelė Teksto analizės įrankių palyginimas

| | SDL MultiTerm Desktop S2011 | SynchroTerm | WordFast | TexNet32 | Projektuojama sistema |
|--|-----------------------------------|-------------|----------|----------|--------------------------|
| Galimybė analizuoti Lietuvišką tekstą | + | - | - | - | + |
| Galimybė pateikti lemuotas frazes | - | + | - | - | + |
| Galimybė iš teksto išrinkti tam tikro ilgio ar struktūros frazes | - | - | - | - | + |
| Galimybė vartotojui sudaryti paieškos tekste šablonus | - | - | - | - | + |
| Sistema gali analizuoti: .doc, .docx, .txt, .pdf failų formatus | ++ | ++ | ++ | ++ | + |

1.5. Frazių ir terminų atpažinimo algoritmai

Per pastaruosius metus, ryškus progresas buvo padarytas, tekstų analizavime norint iš jo išgauti visą reikiamą informaciją, kurią būtų galima panaudoti įvairiems tikslams. Pradedant nuo paprasčiausių frazių išskyrimo algoritmų kurie turi didelį įvairių simbolių ir išsireiškimų skaičių kurie reikalingi identifikuoti ir sugrupuoti tekstą taip, kad suprastų žmogus apie ką šiame tekste yra kalbama neskaičius jo.

Šioje dalyje bus kalbama apie žodynų sudarymo ir teksto supratimo algoritmus. Norint sudaryti specializuotą žodyną reikia iš teksto išskirti raktinius žodžius bei frazes, kurių pagalba bus sudarinėjamas žodynas. Dažniausiai, žodynų įrašų vardai ir aprašai gali programoms padėti sudaryti specializuotus žodynus.

Žodynų įrašai skiriasi nuo techninių vartojamų žodžių. Nes tarkime turėdami vieną žodį paprastame žodyne, tai jo prasmė gali skirtis nuo tarkime techniniame žodyne esančio žodžio.

1.5.1. Variantų atrinkimo algoritmo analizė

Šio algoritmo veikimo principas yra iš teksto išrinkti ir atrinkti panašius žodžiu bei frazes ir juos suskirstyti į tam tikras dalis, iš kurių būtų galima sudaryti specializuotą žodyną. Po to kai žodžiai yra atrenkami, jie sudedami į bendro panašumo logikas, kurios tampa raktiniais žodžiais frazių paieškoje tekste. Toliau apžvelgsime kaip apskaičiuoti galimybę iš gautų žodžių sąrašo ir taip išrinkti frazes. Yra keli tipai formų suskirstyti frazes:

- Simbolių variacijos
- Variantų sudarymas
- Variantai su klaidomis
- Santraukos

Simbolių variacijos algoritmas

Algoritmo logika yra tokia, jog reikia išskirti iš žodžių įvairius skirtukus tarkime: , - , / , ir kt. Kaip žinoma lietuvių kalboje taipogi yra žodžių ar išsireiškimų kurių reikšmė yra viena tačiau jie per vidurį yra išskirti brūkšneliu, tai šis algoritmas tiesiog panaikina tarp žodyje esantį skyriklį ir jį implimentuoja kaip vieną ištisą žodį, tarkime:

„*morfologinė-sintaksinė daryba*“ -> „*Morfologinė sintaksinė daryba*“

Variantų sudarymo algoritmas

Šio algoritmo logika paremta panašių žodžių struktūros ieškojimu. Jis lygina žodžius su kitais tariamais panašiais žodžiais ar frazėmis nors jos gali būti ir nepanašios užrašymu tačiau prasmė gali būti tokia pati, tarkime:

Angl.: Passenger Airbag ir passenger air bag,

Mass Airflow Sensor ir Mass Air Flow Sensor;

Variantų su klaidomis algoritmas

Šis algoritmas atpažįsta nors ir klaidingai (su klaidomis) parašytas frazes ir juos prideda į sąrašą, kuris po to yra pridedamas prie variantų sudarymo algoritmo. Šio algoritmo pagrindinis tikslas yra rasti nors ir galimai klaidingas frazes, tačiau jų nepašalinant. Šis algoritmas tinka ir lietuviško teksto analizei bei frazių atrinkimo.

Santraukos algoritmas

Įvairiuose dokumentuose, straipsniuose, o ypač techninės paskirties dokumentuose, yra naudojama daug sutrumpinimų, dėl to reikia ypač atsargiai į tuos sutrumpinimus žiūrėti bei juos išskirti iš teksto. Šioje dalyje bus naudojamas metodas, kuris yra apibrėžtas šioje dalyje [17], norint atpažinti sutrumpinimus. Pvz.:

Angl.: 4H : four-wheel drive high

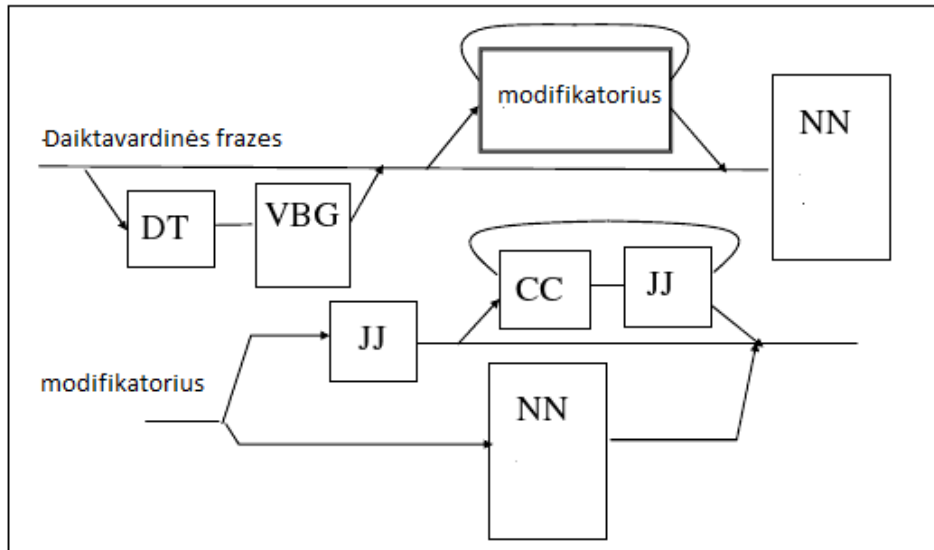
ATF : automatic transmission fluid

Liet.: KTU : Kauno Technologijos Universitetas

UAB : Uždara Akcinė Bendrovė

1.5.2. Daiktavardinių frazių atpažinimo iš sakinio struktūros algoritmas

Žodyno elementas gali būti arba daiktavardinė frazė arba veiksmažodinė frazė. Veiksmažodžius svarstome tik tokius, kurie yra ne pagalbiniai ir imame jų bendrines formas, kurios yra įtraukiamos į žodyną. Daiktavardinėms frazėms yra pristatyta struktūra pagrįsta tyrimais Juston ir Katz(1995). Žemiau pateiktame paveikslėlyje pateikiama kaip atpažįstama daiktavardinės frazės anglų kalboje (9 pav.) [9].



9 Pav. Daiktavardinių frazių struktūra [9]

DT- sprendėjas, VBG- apibūdina specifinę veiksmažodžio formą, JJ- apibūdina būdvardį, CC- apibūdina jungtuką, NN- apibūdina daiktavardžius. Tokia struktūra tinka ir lietuvių kalbai. Pavyzdžiui sakiny "Mes skaičiuojame kalendorines dienas". Frazė bus kalendorinės dienos. Arba " Aš turiu sporto inventoriaus" šiuo atveju frazė būtų sporto inventorių.

1.5.3. Atpažinimas kandidatų į žodyno elementus pagal jų formas

Bendra strategija atpažinti techninę terminiją, tai surasti daiktavardines frazes, o jas rasti galima pagal leksika(mokslinė, buitine), kalbos sudeti(veiksny tarinys papildinys), filtrais. Atpažinimas galimo kandidato į žodyno elementus negalima padaryti tik paprastu leksikos filtru. Taipogi atpažinti žodyno elementą susiduriama su problema dėl dviprasmybių, kuri negali būti išbresta tik leksiniu peržiūrėjimu.

Norint integruoti daugybę paruošiamųjų darbų prieš paleidžiant žodyno gavimą, lingvistikos analize paremta rašyti per tekstą anotacijas(apie žodį trumpa santrauka pvz. kalbos dalis, linksnis). Principe nėra svarbu anotacijų skaičius ar tipas: žodis, vardų sąrašas, specialios frazes, gabaliukai, sakiniai visa tai įeina į anotaciją. Leidžiant per lingvistikos analize, paslepia žodžių išsidėstymo tvarką. Pavyzdys žodyno elementų yra pateiktos žemiau esančioje lentelėje. A-būdvardis, N-daiktavardis ir C-jungtukas.

| Structure | Terms |
|-----------|---|
| AN | genuine part |
| NN | sport utilities |
| AAN | heavy commercial use |
| ANN | rear wiper blade |
| NNN | emission control system |
| AANN | other qualified service technician |
| ACAN | unpaved or dusty roads |
| ANNN | automatic transmission fluid level |
| NNNN | engine oil fluid level |
| AANNN | new personalized oil reset percentage |
| AACAN | certain frontal or near-frontal collision |
| ACAAN | ambient and wide open throttle |
| NNNNN | steering wheel fan speed control |

10 Pav. Kandidatų lentelė

Jeigu žodyno elementas yra naudojamas daug dažniau specifiniame dokumente nei kituose dokumentų kolekcijose, tai tikriausiai rastas elementas yra to dokumento specifinis terminas. Šis algoritmas skaičiuoja tikimybę termino pasirodymo specifiniame dokumente nei paprastame dokumentų tekstuose (*ang.l: text corpus*). Tai užrašoma lygtimi:

$$\frac{\sum_{w_i \in T} \frac{P_d(w_i)}{P_c(w_i)}}{|T|}$$

kur T yra žodžių skaičius elemente, $P_d(w_i)$ tikimybė žodžio w_i specifiniame dokumente, $P_c(w_i)$ tikimybė žodžio w_i dokumentų rinkiniuose. Tikimybės nustatomos pagal dokumentų didžius.

1.5.4. Frazių ir terminų radimo algoritmų analizės išvados

- 1) Atlikta algoritmų analizė parodė jog nei vienas iš analizuotų algoritmų nesuteikia galimybės daiktavardinių bei veiksmažodinių frazių ieškoti lietuviškame tekste.
- 2) Algoritmų analizė jog frazes galima ieškoti naudojant frazių sudedamąsias struktūras.
- 3) Atlikta analizė parodė jog lietuviškam daiktavardinių ir veiksmažodinių frazių paieškai tinkamas algoritmas yra ieškoti pagal žodžių formas.
- 4) Lietuviško teksto analizei, t.y. daiktavardinių bei veiksmažodinių frazių paieškai atlikti, bus remiamasi žodžių formomis grįsto algoritmu. Formos kuriomis bus atpažįstamos frazės bus minimos kaip – šablonai.

1.6. Siekiamo sprendimo apibrėžimas

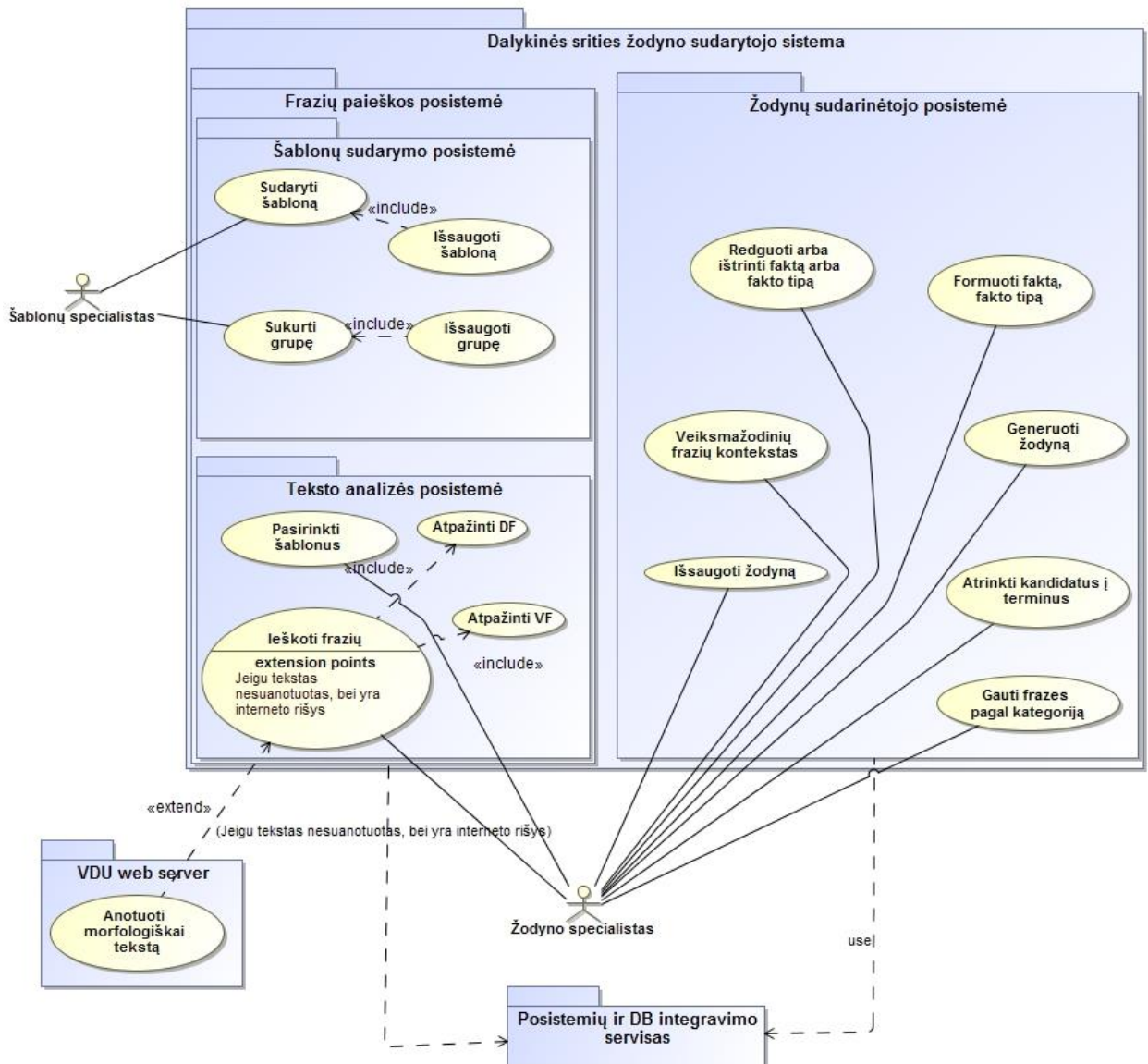
Daiktavardinių bei veiksmažodinių frazių radimo algoritmas, turi pagal vartotojo pateiktus šablonus išskirti iš teksto daiktavardines bei veiksmažodines frazes. Taipogi sistema turėtų teikti galimybę iš teksto išskirti, vardus ir pavardes, sinonimus, sutrumpinimus.

Bendruoju atveju teksto analizės sistema turi atlikti šiuos veiksmus:

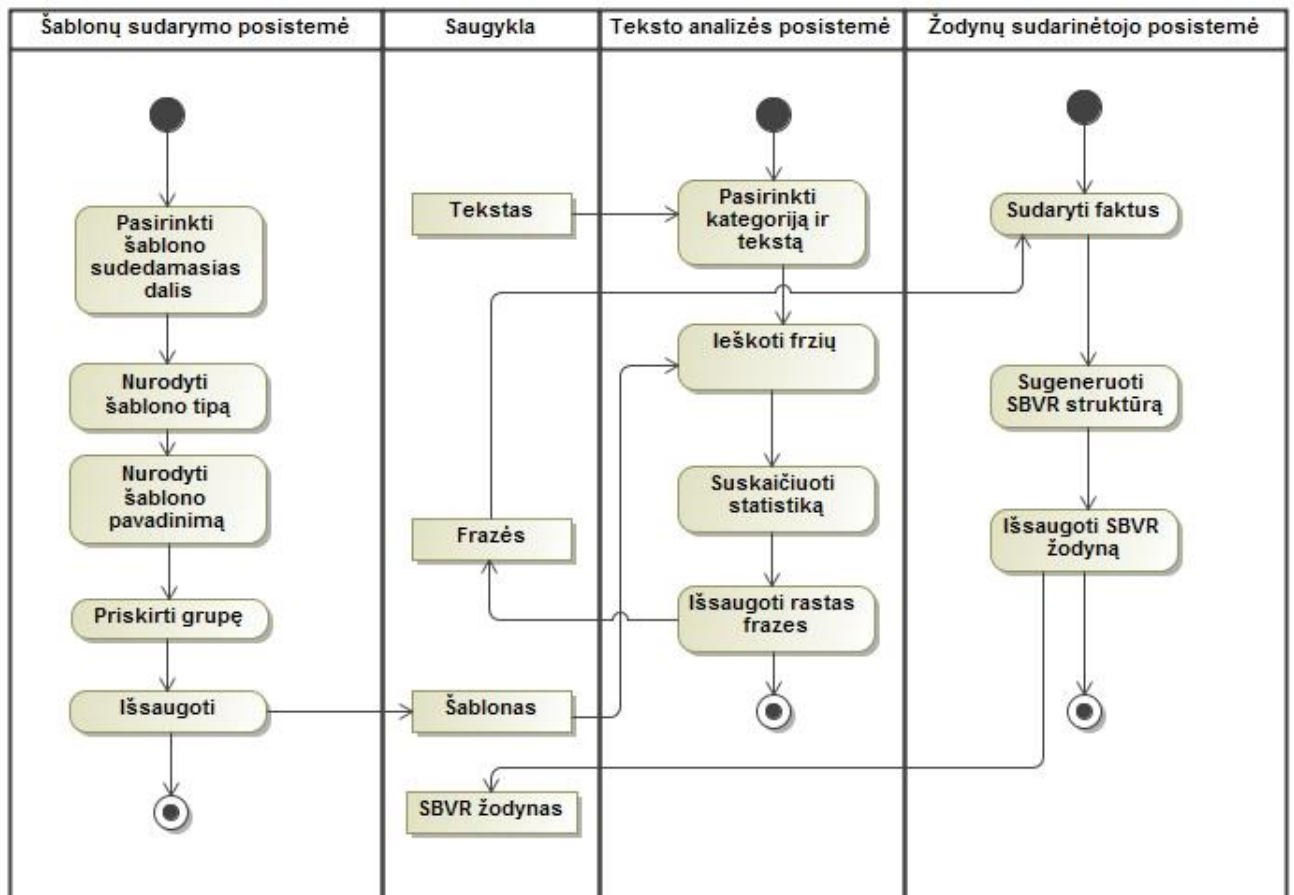
- Leisti vartotojui sudaryti bet kokio tipo šablonus;
- Pripažinti bei ieškoti frazių nepriklausomai nuo šablono sudėtingumo;
- Suteikti galimybę vartotojui peržiūrėti rastus įrašus;
- Leisti įrašus redaguoti bei juos saugoti;

1.6.1. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama

Panaudojimo atvejų diagramoje pavaizduojami pagrindiniai komponentai kurie sudaro dalykinės srities žodyno sudarinėtojo sistemą. Kaip matyti iš (11 Pav. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama) sistemą sudaro dvi pagrindinės dalys, tai yra: Frazijų paieškos posistemė bei žodynų sudarinėtojo posistemė kuri skirta SBVR žodynui generuoti.



11 Pav. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama



12 Pav. Dalykinės srities žodyno sudarytojo sistemos veiklos diagrama

Sistemos veiklos diagramoje yra pavaizduoti pagrindiniai panaudojimo atvejų veiklos etapai (12 Pav. Dalykinės srities žodyno sudarytojo sistemos veiklos diagrama). Sistemos veiklos diagrama išskirta į keturias pagrindines t.y. šablonų sudarinėjimo posistemę, saugyklą, teksto analizės posistemę bei žodynų sudarinėjimo posistemę. „Ieškoti frazių“ panaudojimo atvejų detalesnis pavaizdavimas pateiktas 3 Pav. Frazių atpažinimo veiklos diagrama

1.6.2. Dalykinės srities žodyno sudarinėtojo funkciniai reikalavimai

Pagrindiniai sistemos funkciniai reikalavimai:

Sudaryti šabloną:

Vartotojas turi turėti galimybę sudaryti šablonus, pagal kuriuos tekste ieškos frazių. Taipogi vartotojui sudarant šablonus turi būti galimybė, frazių šablonus išskirti į dvi grupes t.y. *daiktavardinių frazių šablonas* ar *veiksmažodinių frazių šablonas*.

Išsaugoti šabloną:

Vartotojas turi turėti galimybę išsaugoti šablonus. Šablonus vartotojas turi galėti išsaugoti keliose vietose t.y. išsaugoti asmeniniame kompiuteryje, arba jei yra interneto ryšys išsaugoti duomenų bazėje.

Sukurti grupę:

Vartotojas turi turėti galimybę sukurti grupę, į kurią gali pridėti sudarytus šablonus.

Išsaugoti grupę:

Vartotojas turi galėti išsaugoti sudarytas grupes bei išsaugoti priskirtus šablonus prie grupių. Šis veiksmas yra galimas tik tada kai yra interneto ryšys, kadangi visi duomenys apie šablonus bei grupes yra saugomi duomenų bazėje.

Gauti šablonus:

Vartotojas turi turėti galimybę gauti šablonus, kuriuos jis yra sudaręs.

Ieškoti frazių:

Vartotojui turi būti sudaryta galimybė ieškoti frazių tekste. Frazių paieška turi būti išskaidyta į kelias grupes, tokias kaip: Vardų bei pavardžių paieška, sutrumpinimų paieška, sinonimų paieška bei daiktavardinių ir veiksmažodinių frazių paieška naudojant vartotojo sudarytus šablonus.

Ieškoti DF (daiktavardinių frazių):

Sistema turi galėti ieškoti daiktavardinių frazių, naudojant vartotojo sudarytus šablonus.

Ieškoti VF (veiksmažodinių frazių) :

Sistema turi galėti ieškoti veiksmažodinių frazių, naudojant vartotojo sudarytus šablonus.

Anotuoti morfologiškai tekstą:

Sistema turi pati sugebėti išanotuoti pateiktą tekstą. Sistema gali nepriimti, arba neleisti anotuoti teksto dėl interneto ryšio nebuvimo, kadangi dėl teksto anotavimo sistema kreipiasi į nutolusį VDU servisą.

Sudaryti faktą:

Sistema turi galėti bei vartotojui leisti sudaryti faktus, iš rastų daiktavardinių bei veiksmažodinių frazių.

Išsaugoti faktą:

Sistema turi leisti vartotojui išsaugoti sudarytus SBVR faktus, į specialų SBVR duomenų failą, kad sudarytą žodyną būtų galima naudoti ir kitose sistemose.

1.6.3. Dalykinės srities žodyno sudarinėtojo sistemos nefunkciniai reikalavimai**1.6.3.1. Reikalavimai sistemos išvaizdai**

Bendri reikalavimai teminės srities žodyno sudarinėtojo vartotojo sąsajai:

- Lengvai suprantama sąsaja;
- Neįkyri sąsaja (pvz.: nereikalauja pastoviai ką nors vis patvirtinti);
- Aiškus bei patogus pagrindinis valdymo meniu;
- Sąveikaujanti sąsaja;
- Suprantama ir aiškiai išdėlioti bei paaiškinti funkcionalumo mygtukai;

1.6.3.2. Reikalavimai panaudojamumui

Panaudojamumo paprastumas (lengvumas) teminės srities žodyno sudarinėtojo kriterijai:

- Klaidų poveikių mažinimas;
- Galimybė bet kada gauti naujausius sistemos atnaujinimus;
- Procesų vykdymo vizualizacija;
- Galimybė patogiai vartotojui pateikti tiek patį šabloną, tiek rastas frazes;
- Patogus ir lengvai suprantamas šablonų kūrimo sąsaja.
- Patogi ir lengva frazių bei faktų redagavimo sąsaja;

1.6.3.3. Reikalavimai veikimo sąlygoms

Sistema turi veikti windows 32 ir windows 64 bitų operacinėse sistemose: windows xp, windows vista, windows 7, windows 8, windows 8.1.

Norint pilnai naudotis sistema, būtina turėti interneto ryšį. Interneto ryšys būtinas tam, jog sistema galėtų susisiekti su servisu, kuris aptarnauja programą duomenimis iš duomenų bazės.

1.6.3.4. Saugumas

Toliau pateikiami teminės srities žodyno sudarinėtojo sistemos saugumo aspektai:

- Konfidencialumas – sistemoje saugomi tik išanalizuotų tekstų duomenys, nesaugoma jokia kita su asmeniu ir jo tapatybę galinti atskleisti informaciją.
- Privatumas – kiekvienas vartotojas parsijunkdamas prie sistemos identifikuojamas atskirai, tad nėra galimybės pamatyti svetimo vartotojo duomenų.
- Pasiekiamumas – pagrindiniai sistemos duomenys su kuriais operuoja sistemos yra rastos frazės bei tekstai ir vartotojo sudaryti šablonai. Visi šie duomenys dėl patikimumo ir didesnio pasiekiamumo saugomi duomenų bazėje.

1.7. Analizės išvados

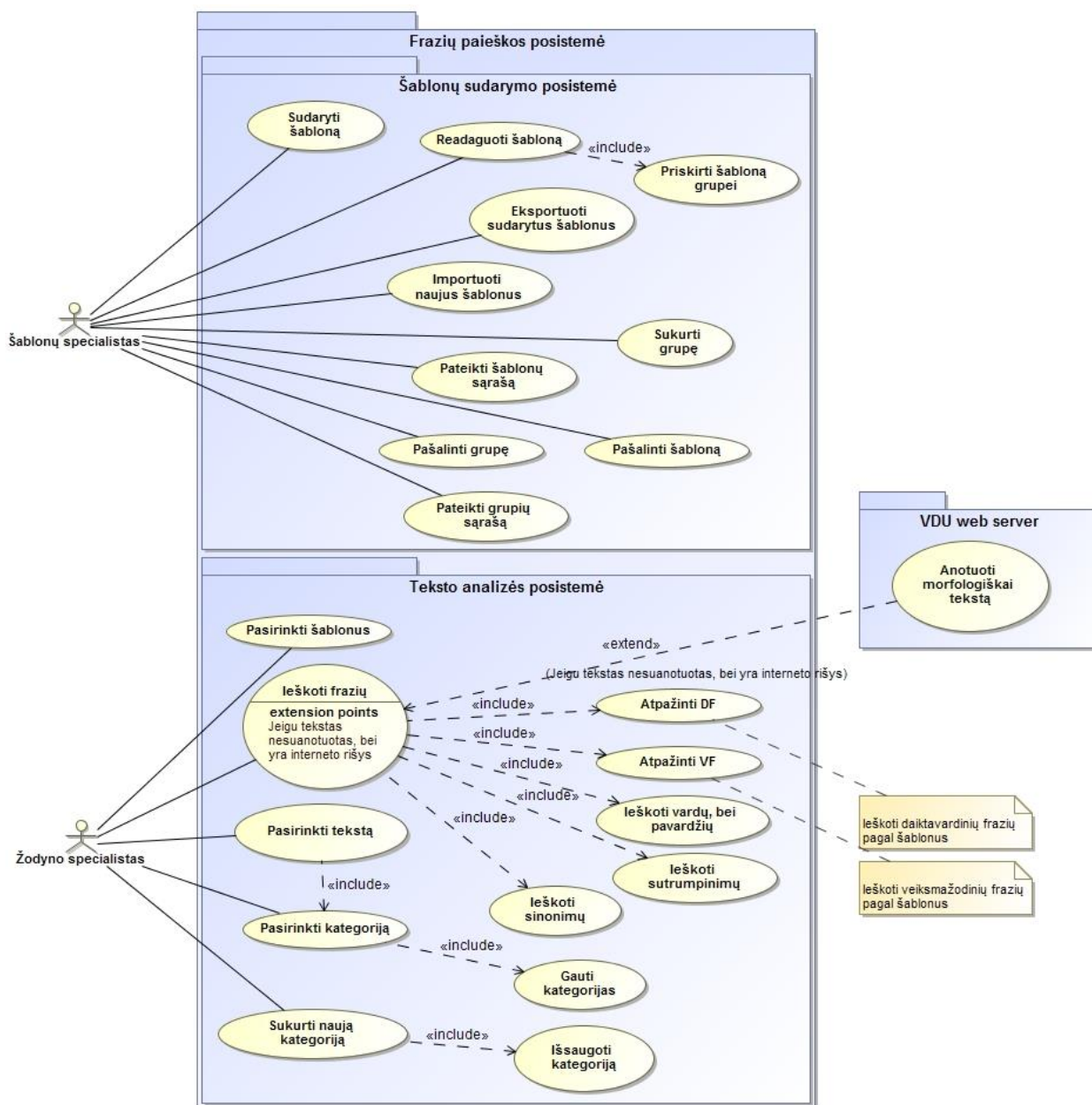
- 1) Atlikta analizė leidžia teigti, jog frazių ar žodžių junginių radimas tekstuose ar žodynų sudarinėjimas yra gan sudėtingas ir laiko reikalaujantis mechanizmas.
- 2) Atlikus panašių sistemų analizę, buvo pastebėta, jog nei viena iš analizuotų sistemų netinka lietuviško teksto analizei. Taipogi nei viena iš sistemų neturėjo galimybės pačiam vartotojui nustatyti ieškomų frazių, junginių ar kitokių struktūrų žodžių, sudaryti numanomus šablonus pagal kuriuos sistema galėtų rasti bei aptikti frazes ar junginius.
- 3) Iš atliktų teminės srities žodynų sudarinėtojų programų analizės galima teigti, jog lietuviško teksto žodynų sudarymas gali būti gan sudėtingas. Norėdami palengvinti frazių radimą bei pačios programos kūrimą galima bus remtis, jau esamomis kitoms kalboms skirtomis funkcijomis bei paieškos metodais.

2. FRAZIŲ PAIEŠKOS POSISTEMĖS SPECIFIKACIJA IR ANALIZĖ

2.1. Reikalavimų specifikacija

Pateiksime (11 Pav. Dalykinės srities žodyno sudarytojo sistemos panaudojimo atvejų diagrama) sistemos panaudojimo atvejų „Frazių paieškos posistemė“ detalesnį vaizdą bei panagrinėsime ir aprašysime pagrindinius funkcinis bei nefunkcinis reikalavimus.

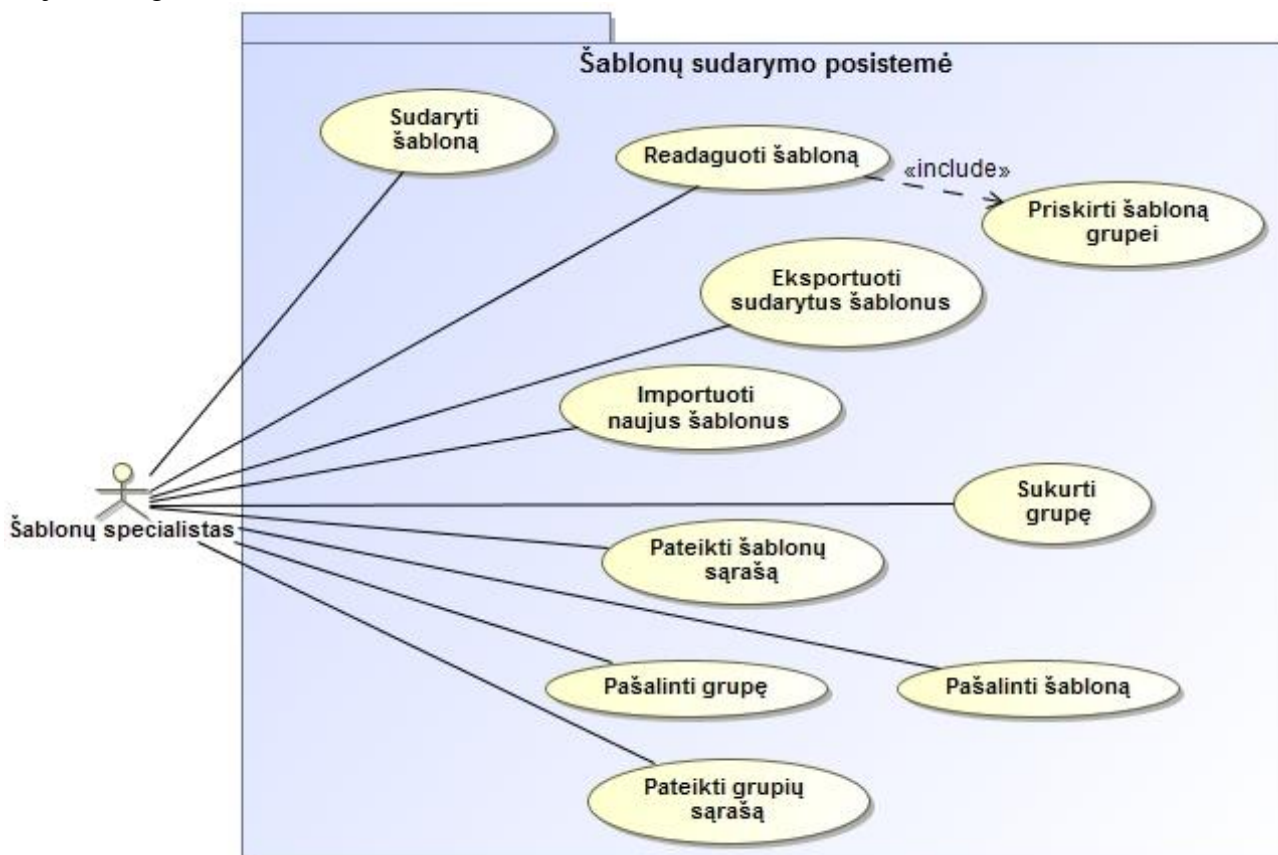
Žemiau (13 Pav. Frazių paieškos posistemės panaudojimo atvejų diagrama) pateiktoje „Frazių paieškos posistemėje“ matyti, jog ji išsiskiria į du paketus. Toliau panagrinėsime abiejų šių paketų funkcinis reikalavimus, nefunkcinis reikalavimus, taipogi sekų ir veiklos diagramas.



13 Pav. Frazių paieškos posistemės panaudojimo atvejų diagrama

2.1.1. Šablonų sudarinėjimo posistemė

Žemiau pateiktoje panaudojimo atvejų diagramoje 14 Pav. Šablonų sudarinėjimo posistemės panaudojimo atvejų diagrama matome „Šablonų sudarinėjimo posistemę“. Toliau bus pateikta pagrindinių panaudojimo atvejų, funkciniai bei nefunkciniai reikalavimai, taipogi jų sekos ir panaudojimo diagramos.



14 Pav. Šablonų sudarinėjimo posistemės panaudojimo atvejų diagrama

2.1.1.1. Šablonų sudarymo posistemės funkcijos

Sudaryti šabloną – vartotojas turi turėti galimybę sudaryti, bet kokio tipo bei sudedamųjų dalių šabloną. Šablono sudarymo metu vartotojas turi galėti, priskirti šablonui vardą, nurodyti jo tipą. Jei sistema yra prijungta prie interneto ir yra pasiekiamas servisas, vartotojas tada turi turėti galimybę sudarytus šablonus išsaugoti duomenų bazėje.

Redaguoti šabloną – vartotojui turi būti suteikta galimybė redaguoti pasirinktą šabloną. Vartotojas pasirinkęs šabloną turi galėti, keisti šablono pavadinimą, grupę kuriai priklauso ir šablono tipą.

Priskirti šabloną grupei – vartotojas turi turėti galimybę priskirti šabloną grupei. Grupės tikslas yra atskirti tam tikrus frazių paieškos šablonus, pagal vartotojo reikalavimus. Taipogi vartotojui turi būti galima tam tikras jau sudarytas frazes priskirti ir kitoms grupėms, ar šablonus pašalinti iš kitų grupių.

Eksportuoti sudarytus šablonus – galimybė išeksportuoti pasirinktus šablonus į XML formatą.

Importuoti sudarytus šablonus – galimybė importuoti XML formate esančius šablonus.

Sukurti grupę – vartotojas turi turėti galimybę sukurti naują grupę. Jeigu sistema turi prieigą prie interneto ir taipogi servisas funkcionuoja, tai vartotojas turi galėti naujas suskurtas grupes išsaugoti duomenų bazėje.

Išsaugoti šabloną – sistema turi galėti išsaugoti vartotojo sudarytus šablonus.

Gauti grupes – sistema turi vartotojui pateikti, būtent jo sudarytas grupes bei prie grupių automatiškai priskirti šablonus.

Gauti šablonus – sistema turi vartotojui pateikti visus jo sudarytus šablonus.

Išsaugoti grupę – sistema turi galėti išsaugoti vartotojo sukurtą naują grupę.

2.1.1.2. Šablonų sudarymo posistemės nefunkciniai reikalavimai

- **Reikalavimai išvaizdai:**

Šablonų sudarinėjimo langas turi būti lengvai suprantamas bei aiškus. Vartotojui turi būti aiškiai ir suprantamai išdėstyti elementų t.y. kalbos dalių, giminių, laipsnių ir t.t. pasirinkimas bei pridėjimas į šabloną.

- **Reikalavimai saugumui:**

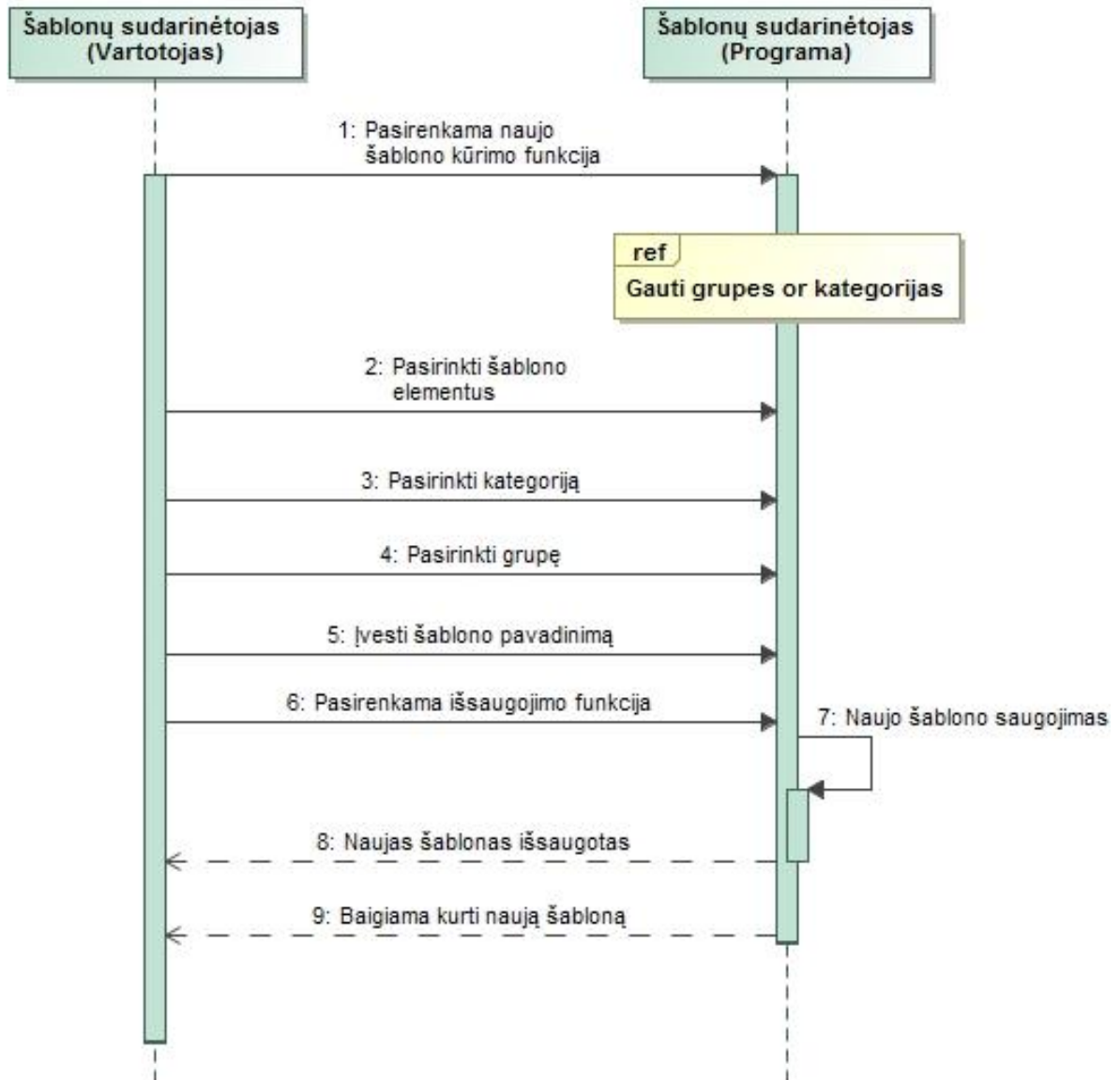
Vartotojai negali matyti vienas kito šablonų, nebent jie juos apsikeičia *eksportuodami* ir *importuodami* kitam asmeniui.

Išeksportuoti frazių šablonai privalo turėti unikalų saugumo raktą, kuris yra tinkamas tik vienam vartotojui, taip apsaugant dėl duomenų nutekėjimo kitam vartotojui.

Duomenų perdavimas iš vartotojo sistemos į duomenų saugojimo saugyklą, privalo būti saugus, kad nebūtų prarandama duomenų.

2.1.1.3. Šablonų sudarymo posistemės sekų diagramos

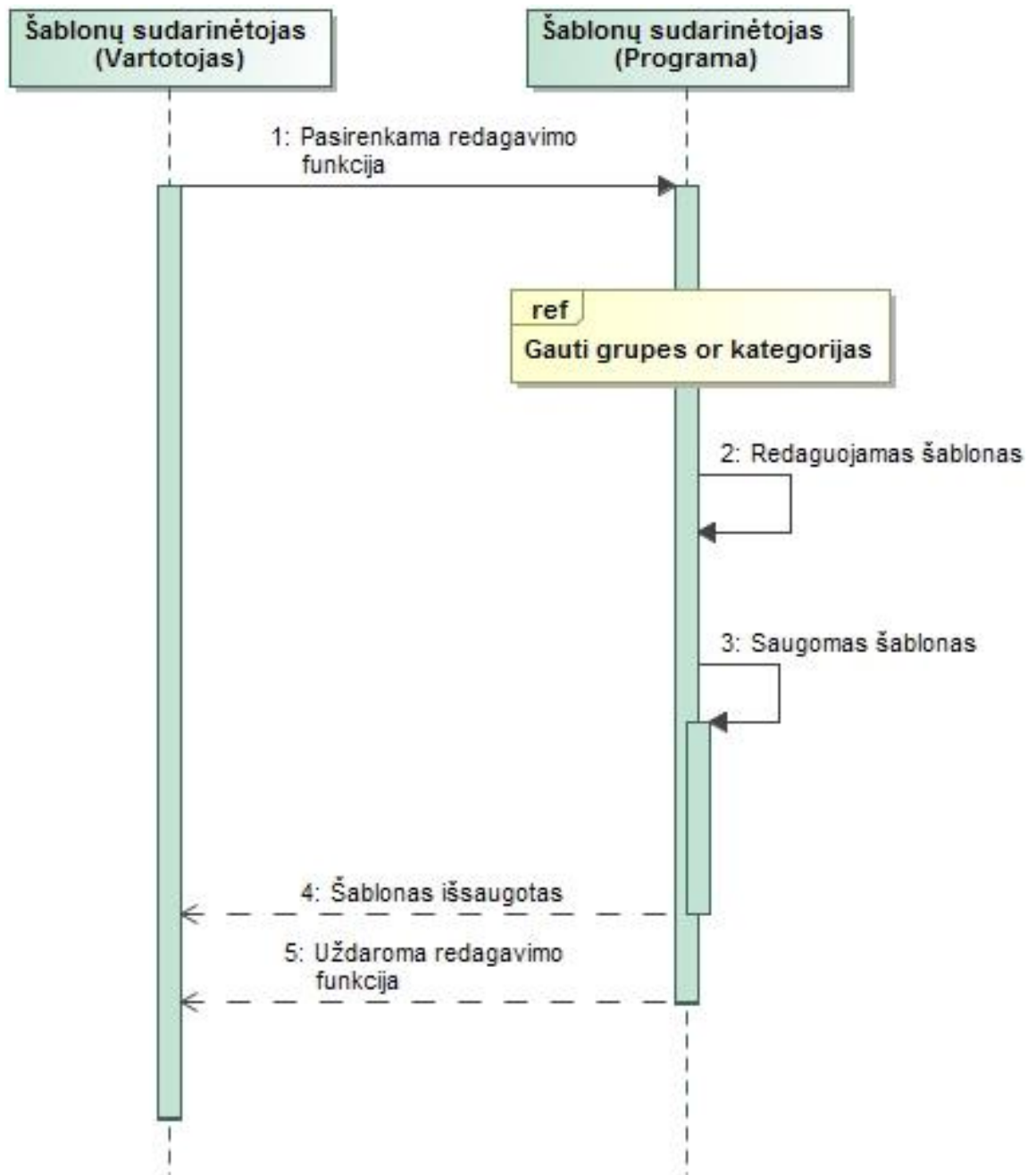
Žemiau esančioje sekų diagramoje detaliau panagrinėsime panaudojimo atvejo „Sudaryti šabloną“ scenarijų, kurio sekų diagrama pateikta (15 Pav. P.A. „Sudaryti šabloną“ sekų diagrama) paveiksle.



15 Pav. P.A. „Sudaryti šabloną“ sekų diagrama

Pagrindiniame sistemos lange pasirinkus funkciją frazių šablono sudarinėtojas yra atidaromas šablonų sudarinėtojo langas. Vartotojas šiame lange gali pasirinkti šablono elementus, t.y. kalbos dalį, giminę, skaičių, linksnį, taipogi pasirinkti kategoriją, įvesti šablono grupę, įvesti šablono pavadinimą. Po to gali visus šiuos pakeitimus išsaugoti. Kada naujas šablonas yra išsaugomas sistema automatiškai išvalo visus sistemos laukus, taip paruošdama langą naujo šablono kūrimui.

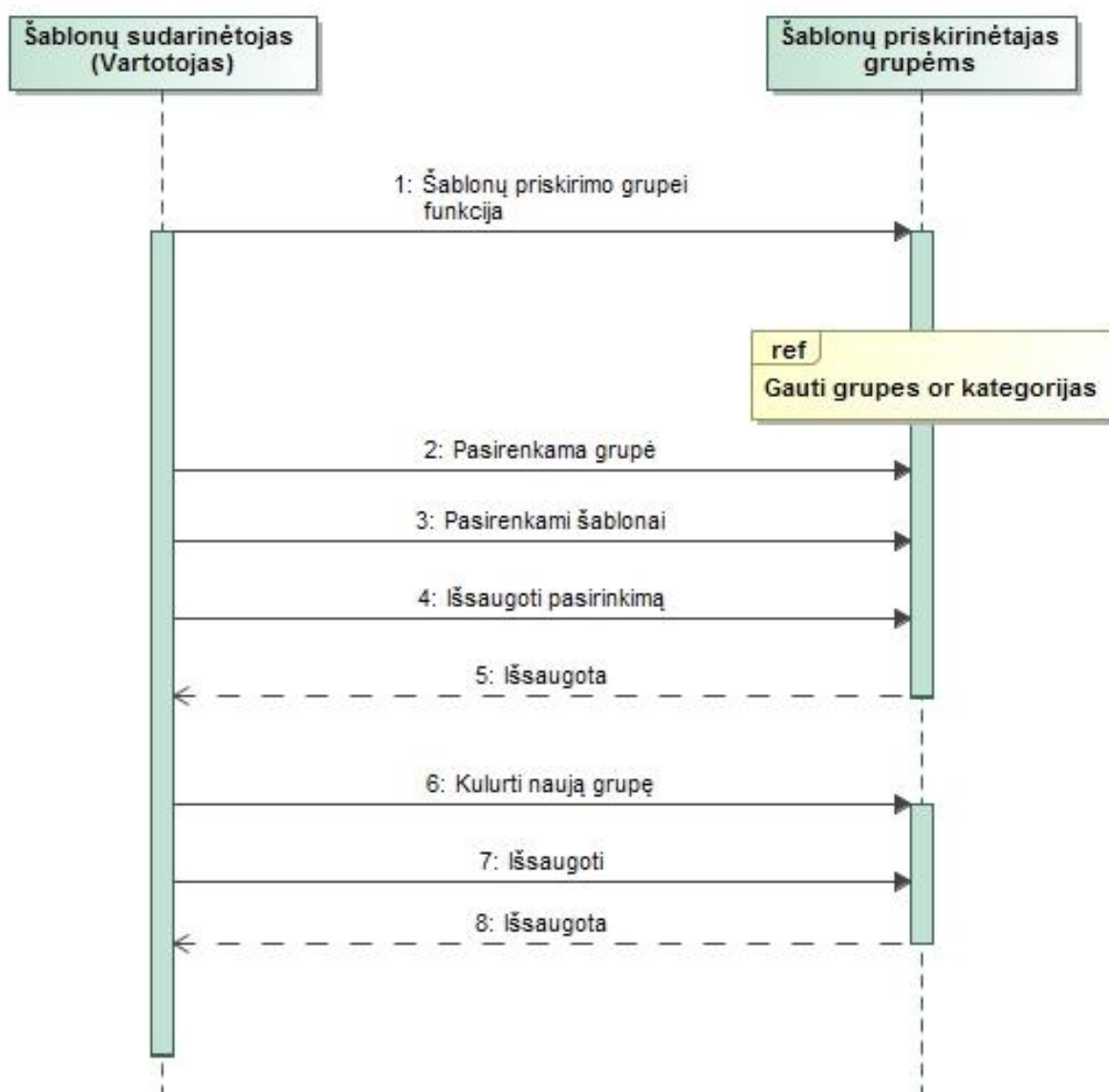
Žemiau esančiame paveiksle detaliau panagrinėsime panaudojimo atvejo „Redaguoti šabloną“ scenarijų, kurio sekų diagrama pateikta (16 Pav. P.A. „Redaguoti šabloną“ sekų diagrama) paveiksle.



16 Pav. P.A. „Redaguoti šabloną“ sekų diagrama

Šablonų sudarinėjimo lange pasirinkus redagavimo funkciją, t.y. pasirinkus tam tikrą šabloną ir nuspaudus redagavimo funkciją, vartotojui suteikiama galimybė redaguoti šabloną. Vartotas redaguojamam šablonui galima pakeisti kalbos dalį, giminę, skaičių ir kita. Taipogi suteikiama galimybė redaguoti pavadinimą, pakeisti grupę, ar šablono tipą. Vartotojui baigus redagavimą, sistema naujus pakeitimus išsaugo duomenų bazėje. Kai pakeitimai išsaugomi, sistema automatiškai užbaigia redagavimą, taip užbaigdama redagavimo funkciją.

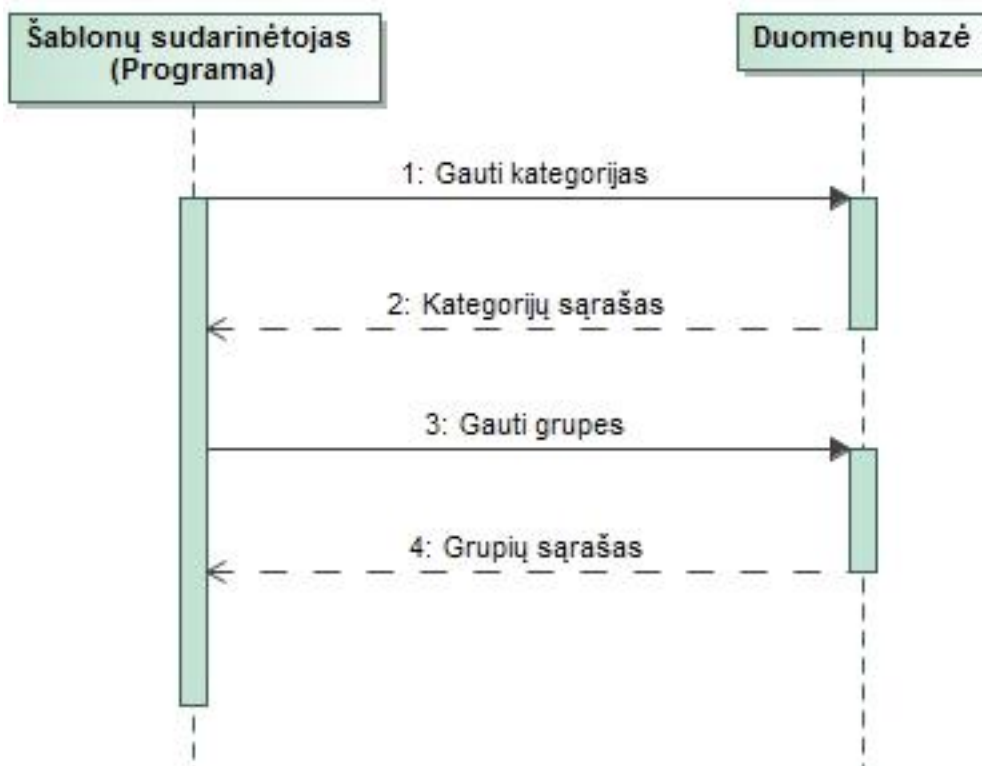
Detalesnė panaudojimo atvejo „Priskirti šablona grupei“ sekų diagrama pateikta (17 Pav. P.A. „Priskirti šablona grupei“ sekų diagrama) paveiksle.



17 Pav. P.A. „Priskirti šablona grupei“ sekų diagrama

Šablonų sudarinėtojas taipogi turi galimybę šablonus priskirti tam tikroms grupėms taip juos sugrupuodamas. Pasirinkus šablonų priskyrimo funkciją vartotojui atsidarymo papildomas langas su vartotojo grupėmis ir šablonais. Vartotojas pirmiausia pasirenką grupę po to pasirenka šablonus kurios nori priskirti būtent šiai grupei. Pasirinkus atitinkamus šablonus specialiai grupei, galima šiuos duomenis išsaugoti. Dar viena papildomą funkcija šablonų priskyrimo funkcijoje yra naujo šablono sukūrimas.

Panaudojimo atvejo „Gauti kategorijas ir grupes“ sekų diagrama pateikiama (18 Pav. P.A. „Gauti kategorijas ir grupes“ sekų diagrama) paveiksle.

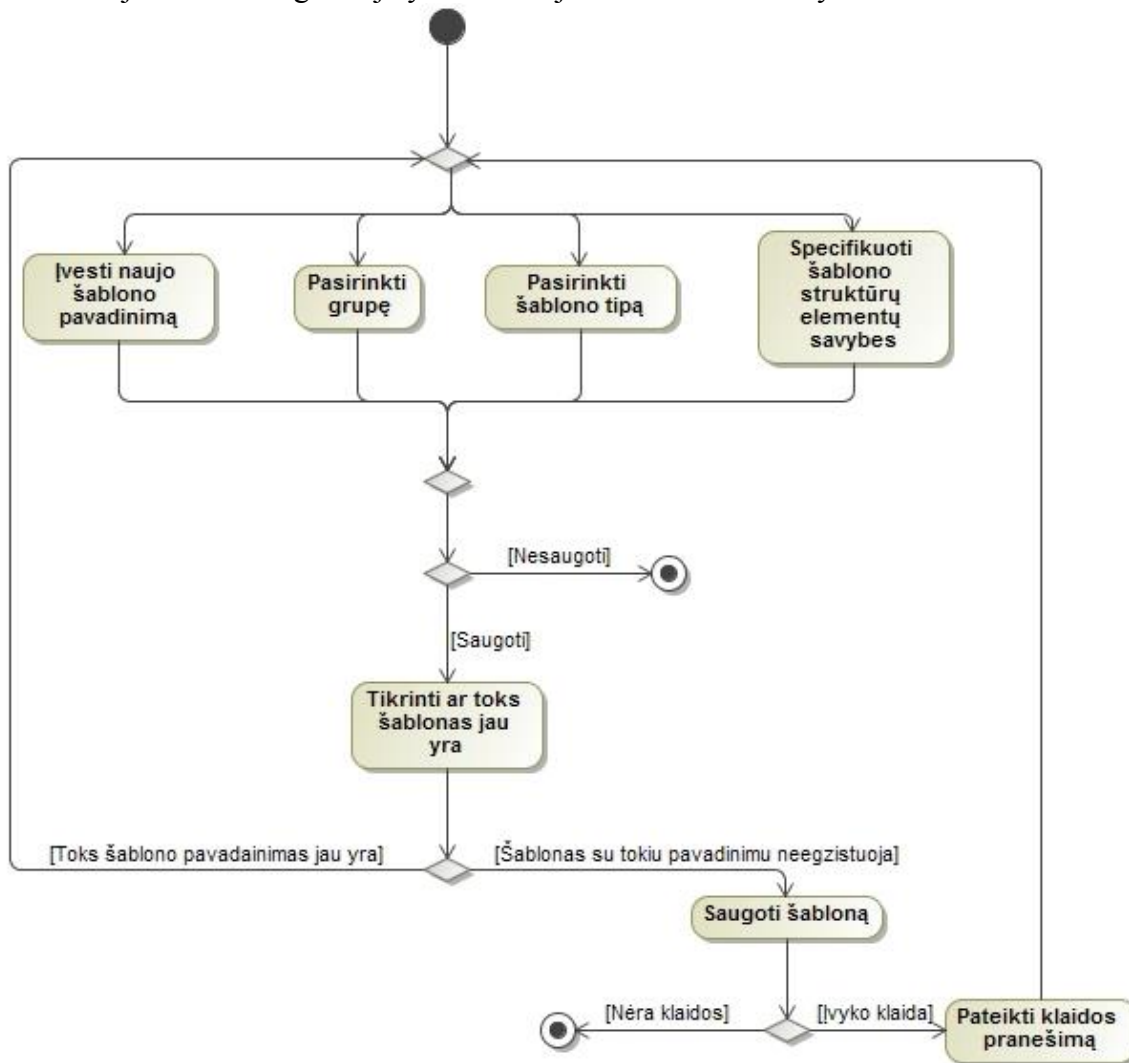


18 Pav. P.A. „Gauti kategorijas ir grupes“ sekų diagrama

Kiekvienas vartotojas turi savo sudarytas, grupes, kategorijas ir taipogi šablonus. Aukščiau pavaizduotoje veiklos diagramoje vaizduojama sekų diagrama kuri parodo kaip sistema gauna kategorijas bei grupes.

2.1.1.4. Šablonų sudarymo posistemės veiklos diagramos

Žemiau esančioje veiklos diagramoje yra vaizduojamas šablono sudarymas:

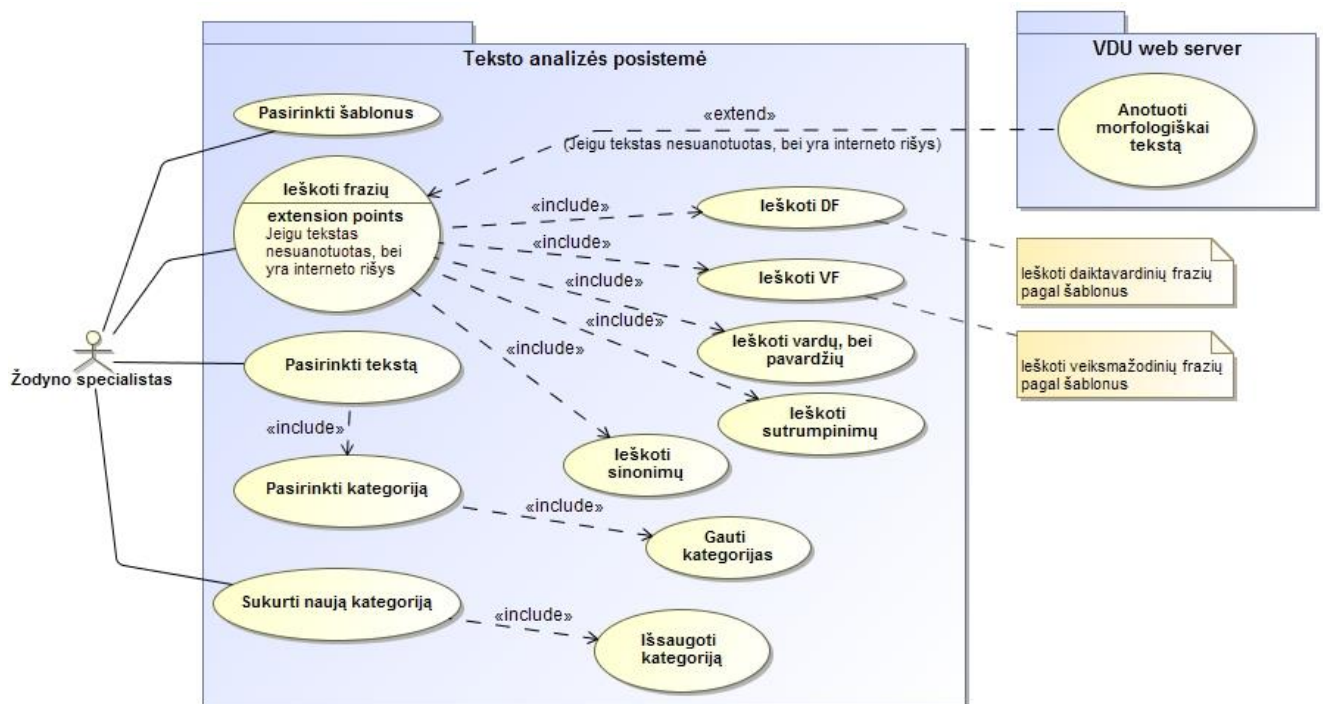


19 Pav. P.A. „Sudaryti šablona“ veiklos diagrama

Šablonas susideda iš keturių pagrindinių dalių, t.y. pavadinimo, grupės, tipo bei struktūrų elementų savybių. Šios keturios dalys yra privalomos kuriant šablonus. Vartotojas pasirinkęs visus elementus gali jų nesaugoti, tačiau jei pasirenkama saugoti, tai sistema automatiškai ar toks šablono pavadinimas jau egzistuoja pas vartotoją. Jei šablonas su tokiu pavadinimu jau egzistuoja tai sistema jį informuoja, jog šablonas su tokiu pavadinimu jau egzistuoja. Jeigu šablonas su tokiu pavadinimu dar neegzistuoja pas vartotoją, tai sistema jį išsaugo, tačiau jei įvyksta sistemoje klaida, šablonas neišsaugomas.

2.1.2. Teksto analizės posistemė

Žemiau pateiktoje panaudojimo atvejų diagramoje (20 Pav. Teksto analizės posistemės panaudojimo atvejų diagrama) matome „Teksto analizės posistemę“. Toliau bus pateikta pagrindinių panaudojimo atvejų, funkciniai bei nefunkciniai reikalavimai, taipogi jų sekos ir panaudojimo diagramos.



20 Pav. Teksto analizės posistemės panaudojimo atvejų diagrama

2.1.2.1. Teksto analizės posistemės funkcijos

Pasirinkti šablonus – vartotojui turi būti suteikta galimybė bei įrankis kuriuo jis galėtų sudaryti šablonus. Vartotojui turi būti pilna galimybė sudaryti tam tikros sudėties šablonus su tam tikrais šablono elementais.

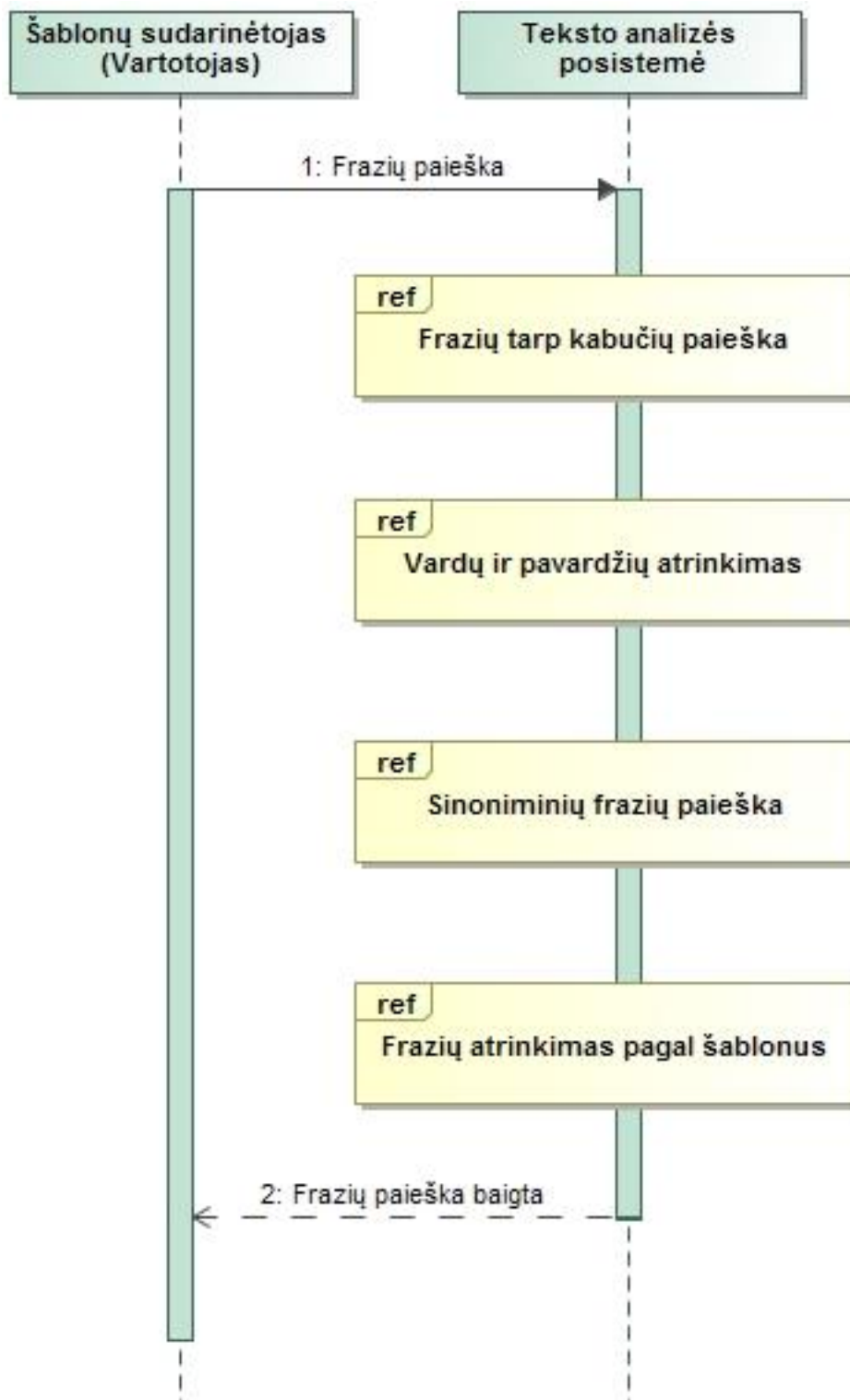
Ieškoti frazių – vartotojui turi būti suteikta galimybė ieškoti įvairaus tipo frazių tokių kaip: *daiktavardinės frazės, veiksmažodinės frazės, vardai bei pavardės, sutrumpinimai* taipogi *sinoniminių* žodžių.

Pasirinkti tekstą – vartotojui turi būti suteikta galimybė įkelti tekstą, ar tekstus kuriuos nori analizuoti. Taipogi sistema turi priimti įvairaus tipo dokumentus: *pdf., doc., txt.* Prieš atliekant frazių paiešką vartotojui yra suteikiama galimybė, savo analizuojamus tekstus ar tekstą priskirti kategorijai.

Sukurti naują grupę – vartotojui turi būti suteikta galimybė susikurti naują grupę, į kurią vartotojas galės priskirti tam tikrus šablonus.

2.1.2.2. Teksto analizės posistemės sekų diagramos

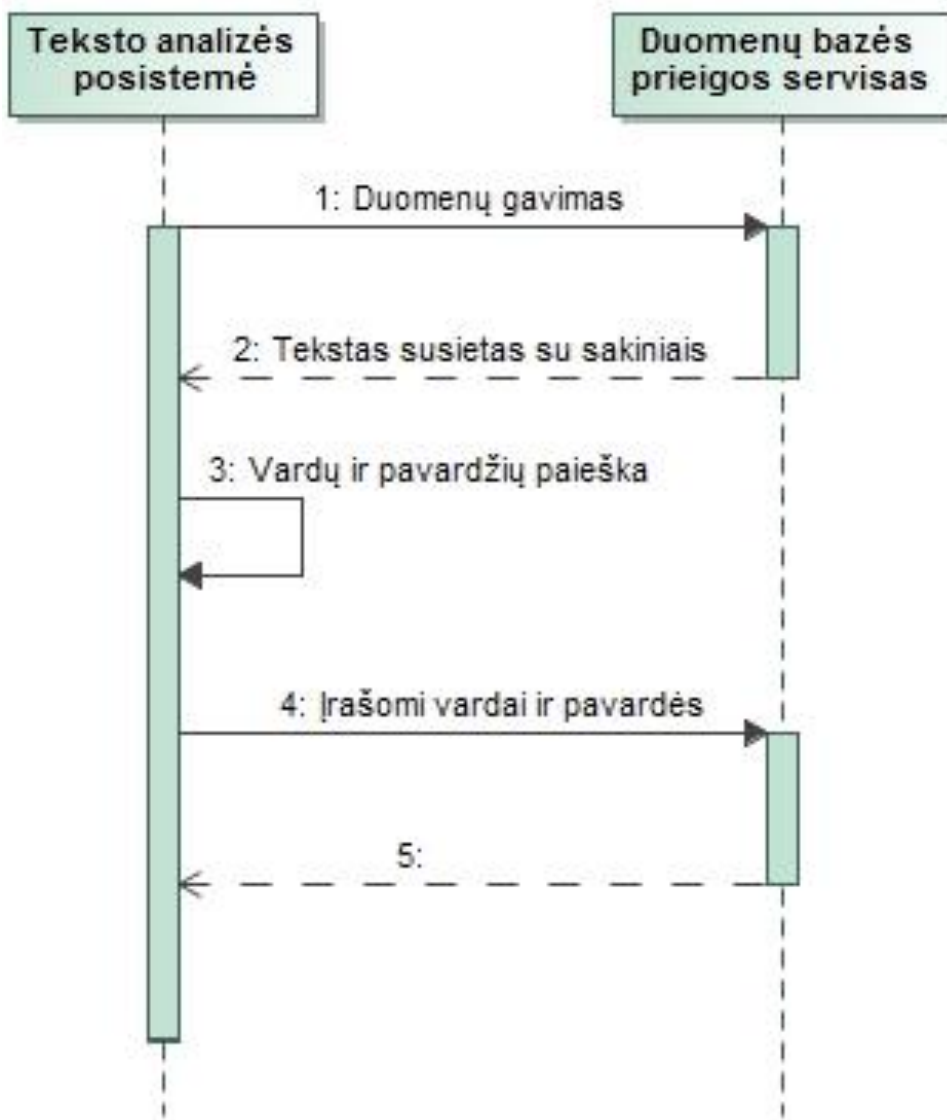
Žemiau esančiame paveiksle (21 Pav. P.A. „Ieškoti frazių“ sekų diagrama) yra vaizduojamas frazių paieškos sekų diagrama, kuri apjungia šiuos pabaudos atvejus: „Frazių tarp kabučių paieška“, „Vardų pavardžių atrinkimas“, „Sinoniminių frazių paieška“, „Frazių atrinkimas pagal šablonus“.



21 Pav. P.A. „Ieškoti frazių“ sekų diagrama

Vartotojas pasirinkęs frazių atrinkimo funkciją, pradeda frazių paiešką. Šioje funkcijoje nuosekliai yra vykdomos 4 funkcijos susijusios su frazių radimu, kurių sekų diagramos yra pateiktos žemiau.

Panaudos atvejo „Ieškoti vardų bei pavardžių“ sekų diagrama pateikta (22 Pav. P.A. „Ieškoti vardų bei pavardžių“ sekų diagrama) paveiksle.

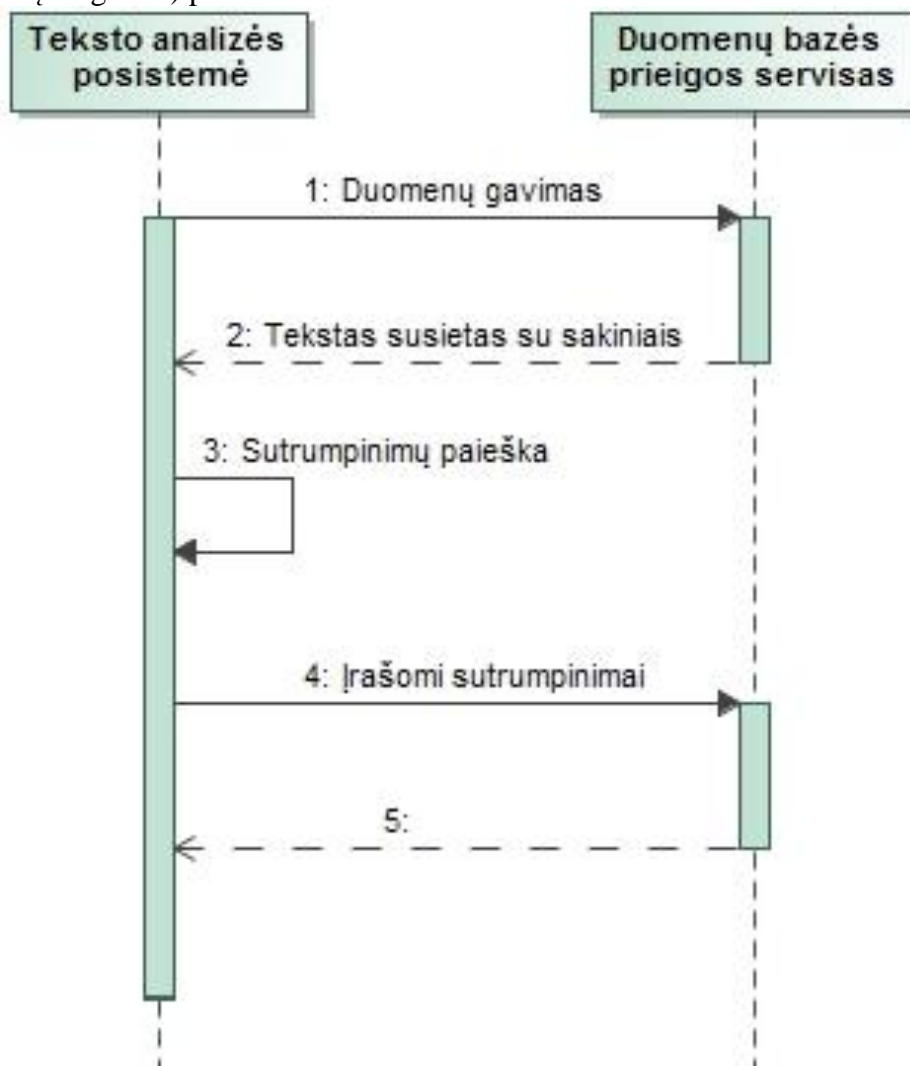


22 Pav. P.A. „Ieškoti vardų bei pavardžių“ sekų diagrama

„Ieškoti vardų bei pavardžių“ pirmiausia iš duomenų bazės atsisiunčia duomenis kuriose yra tekstas bei žodžiai susieti su sakiniais. Toliau vyksta paieška vardų ir pavardžių.

Baigus paiešką visi rasti vardai bei pavardės išsaugomos duomenų bazėje.

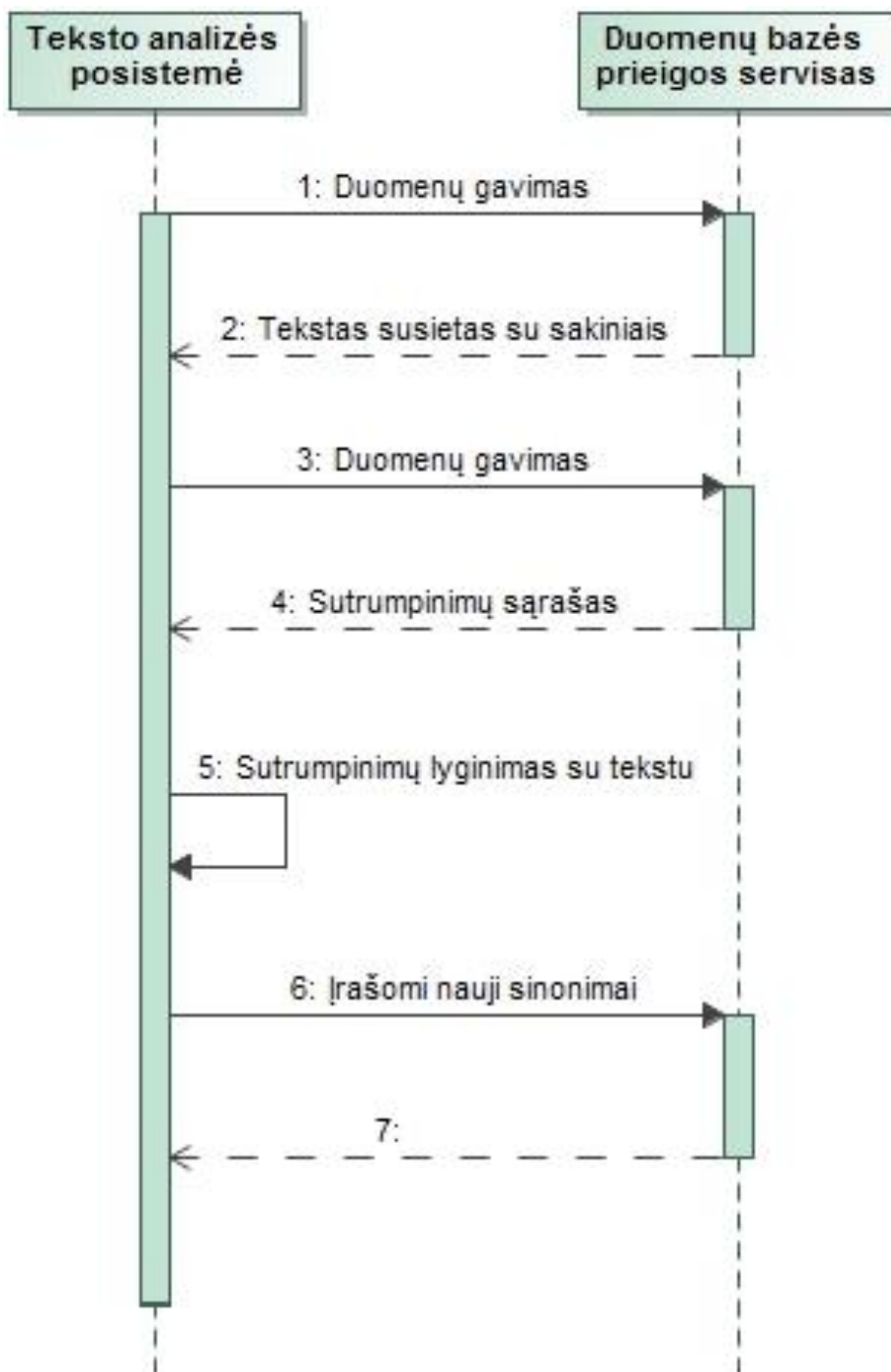
Panaudos atvejo „Ieškoti sutrumpinimų“ sekų diagrama pateikta (23 Pav. P.A. „Ieškoti sutrumpinimų“ sekų diagrama) paveiksle.



23 Pav. P.A. „Ieškoti sutrumpinimų“ sekų diagrama

„Ieškoti sutrumpinimų“ pirmiausia iš duomenų bazės atsiunčia duomenis kuriose yra tekstas bei žodžiai susieti su sakiniais. Po to vyksta sutrumpinimų (KTU, VDU ir t.t.) paieška pagal joje esančias taisykles iki tol kol baigiasi duomenys, kurios atsiuntė iš loklios duomenų bazės. Baigus paiešką visi rasti sutrumpinimai išsaugomi duomenų bazėje.

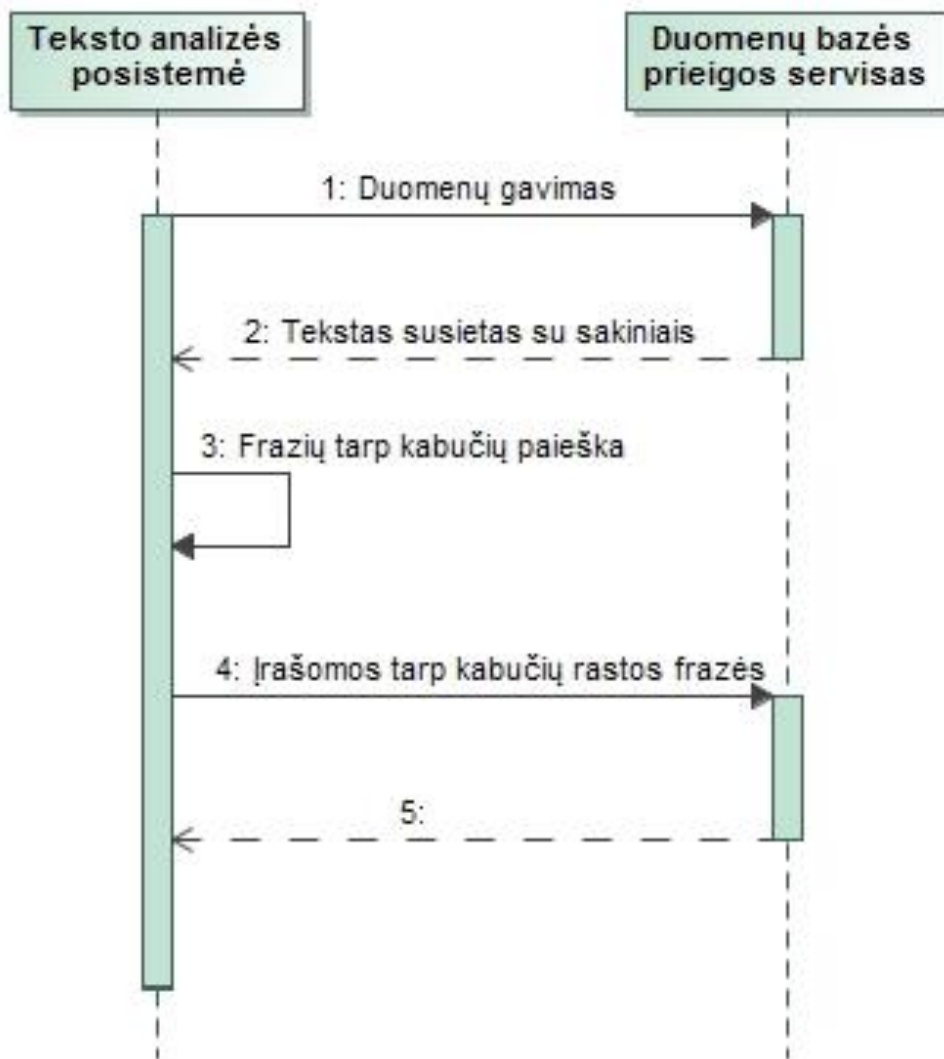
Panaudos atvejo „Ieškoti sinonimų“ sekų diagrama pateikiama (24 Pav. P.A. „Ieškoti sinonimų“ sekų diagrama) paveiksle.



24 Pav. P.A. „Ieškoti sinonimų“ sekų diagrama

Pirmiausia yra gaunamas, tekstas susietas su sakiniais toliau gaunamas sutrumpinimų sąrašas, tada lyginama sutrumpinimai su tekstu, bigę lyginti įrašome sinonimus į duomenų bazę.

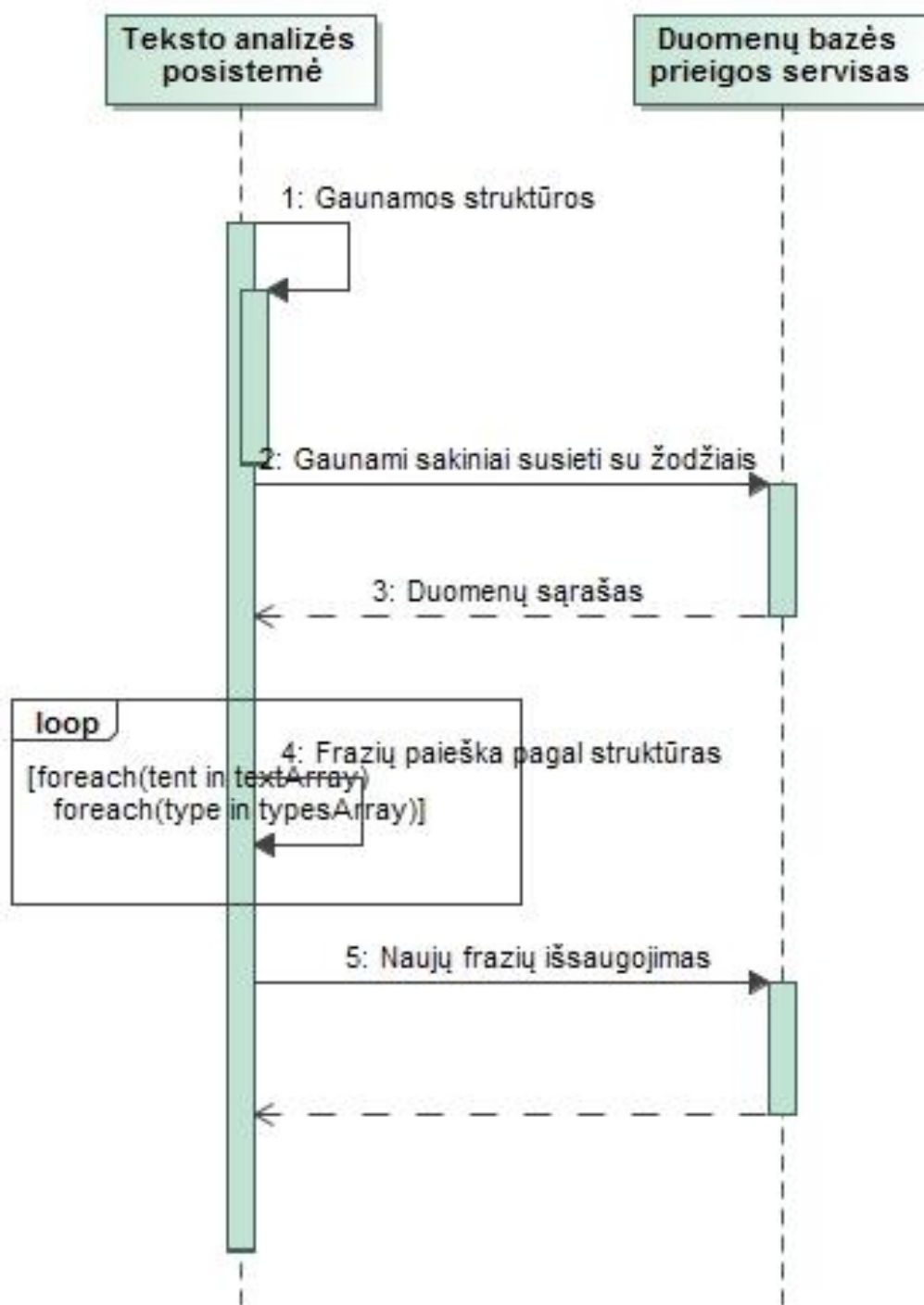
Panaudos atvejo „Ieškoti frazių tarp kabučių“ sekų diagrama pateikiama (25 Pav. P.A. „Ieškoti frazių tarp kabučių“ sekų diagrama) paveiksle.



25 Pav. P.A. „Ieškoti frazių tarp kabučių“ sekų diagrama

Pirmiausia yra gaunamas, tekstas susietas su sakiniiais, toliau ieškome frazių tarp kabučių. Pabaigus ieškoti įrašomos tarp kabučių rastos frazės į duomenų bazę.

Detaliau panagrinėta frazių paieškos sekų diagrama pateikta (26 Pav. „Frazių paieška pagal šablonus“ sekų diagrama) paveiksle. Joje pavaizduota pagrindinė frazių paieška pagal šablonus seka, tarp „Teksto analizės posistemės“ ir „Duomenų bazės prieigos serviso“.

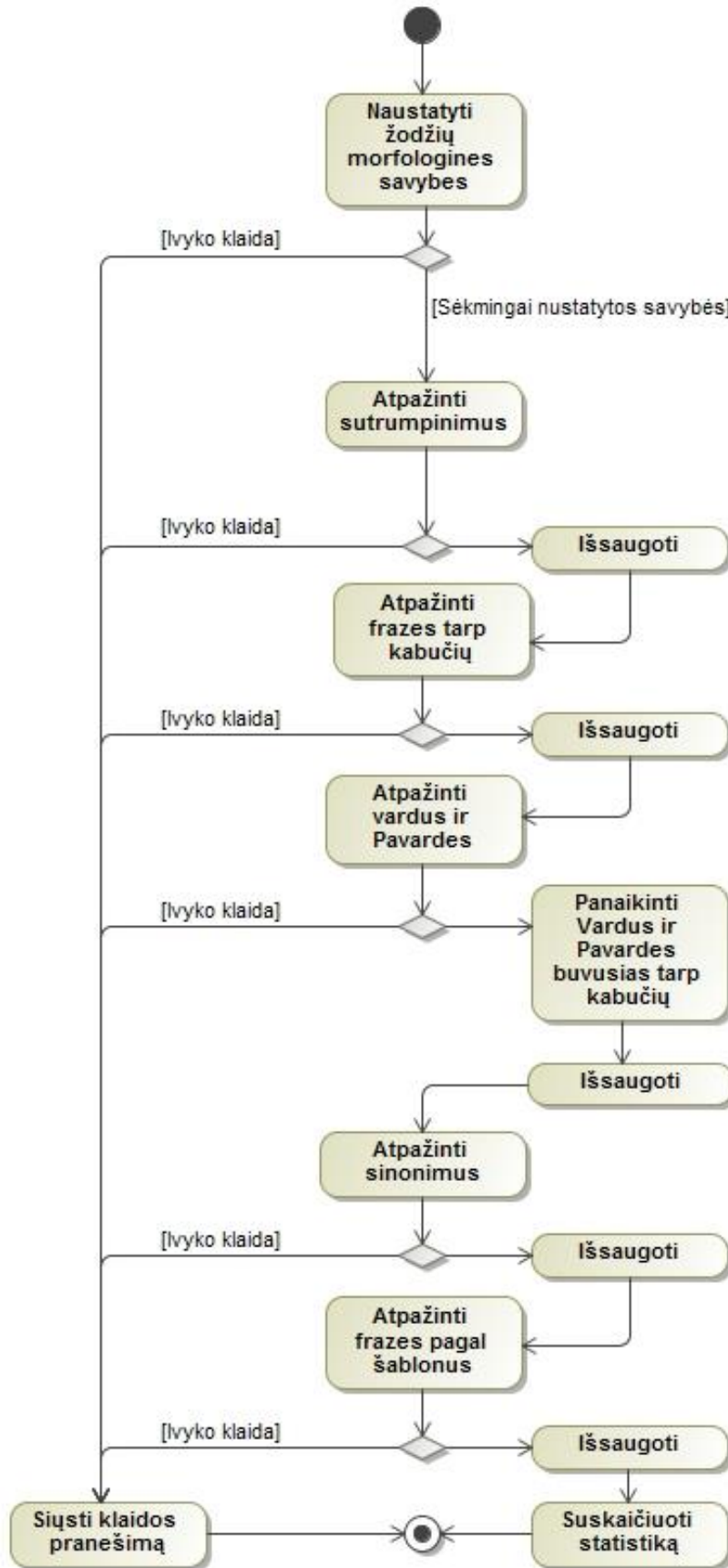


26 Pav. „Frazių paieška pagal šablonus“ sekų diagrama

Frazių paieškos pagal šablonus funkcija pirmiausia iš frazių struktūrų duomenų failo pasiima visas rastas struktūras. Vėliau einama prie sekančio žingsnio „Naujų struktūrų sudarymas“ kuriame neišbaigtos frazių struktūros yra papildomomis, kad užpildytu visą pateiktos struktūros aibę. Vėliau gaunamos iš duomenų bazės duomenys t.y. sakiniai susieti su žodžiais, kurie vėliau yra analizuojami. Prieš paskutinis žingsnis frazių ieškojime tai pagal pateiktus šablonus, tekste ieškoti frazių atitikmenų. Ciklas vyksta tol kol ne iki galo yra išanalizuotas tekstas. Vėliau visos rastos frazės išsaugomos duomenų bazėje.

2.1.2.3. Teksto analizės posistemės veiklos diagramos

Žemiau esančioje veiklos diagramoje yra vaizduojama frazių radimas:

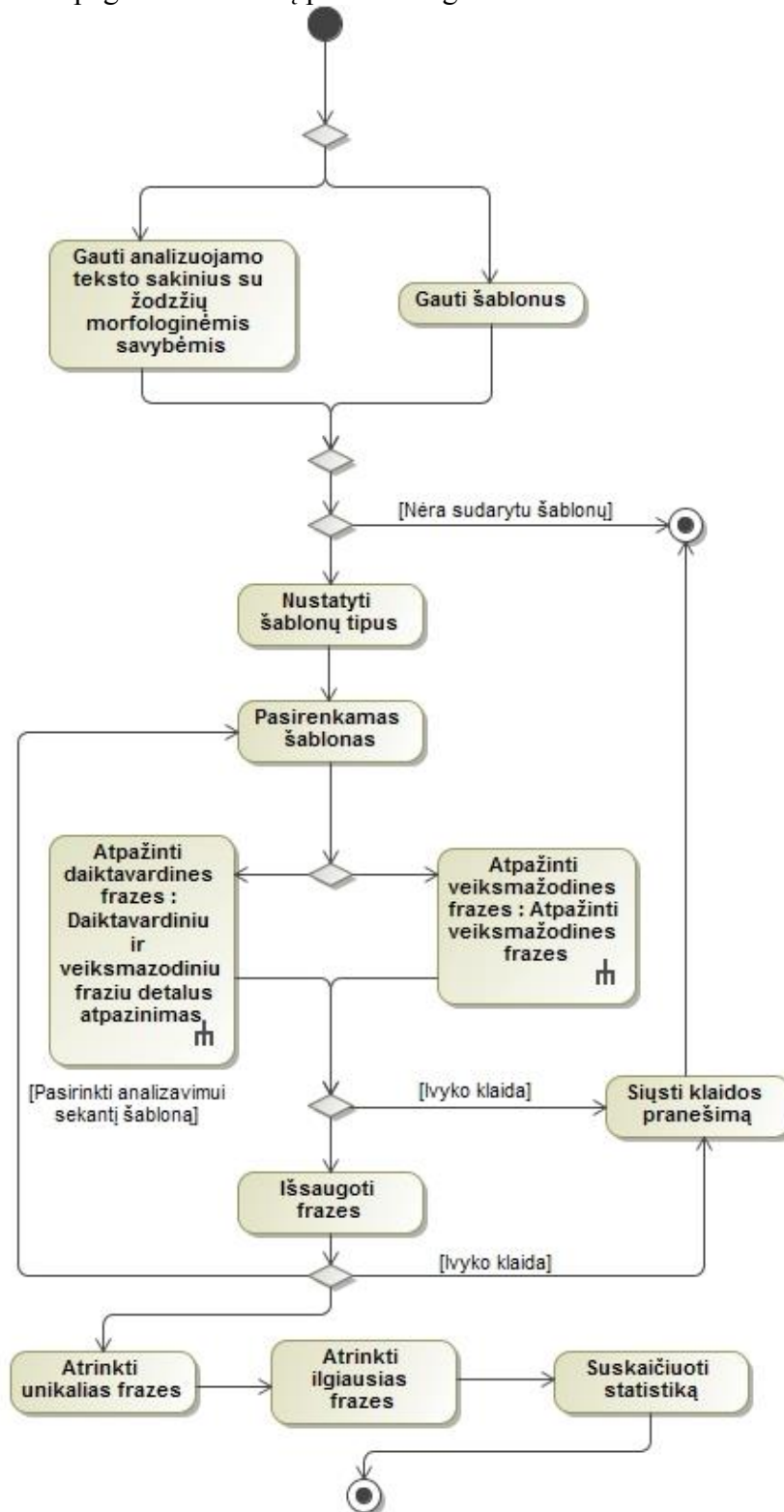


27 Pav. „Frazių radimo algoritmo“ veiklos diagrama

3. AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE SUDARYTŲ ŠABLONŲ PAGRINDU SPRENDIMO ALGORITMAS

3.1. Daiktavardinių bei veiksmažodinių frazių paieškos veikimo aprašymas

Žemiau pateiktame (28 Pav. „Frazių radimo naudojant šablonus algoritmas“ veiklos diagrama) paveiksle, kuriame matome pagrindines frazių paieškos algoritmo veiksmus.



28 Pav. „Frazių radimo naudojant šablonus algoritmas“ veiklos diagrama

Veiklos diagramoje (28 Pav. „Frazių radimo naudojant šablonus algoritmas“ veiklos diagrama) pavaizduota esminė frazių paieškos dalis, kuri yra susieta su šablonų gavimu bei jau rezultatų išsaugojimu. Algoritmas prieš pradėdamas teksto analizę pirmiausias susirenka vartotojo sudarytus bei nurodytus šablonus ir taipogi gauna analizuojamo teksto žodžius su jų morfologinėmis savybėmis. Jei vartotojas nebuvo sudaręs jokių šablonų tai algoritmas nieko neieško ir baigia paiešką. Jei šablonų buvo tai sistema automatiškai nustato jų ilgį pagal kurį po to atliekama analizė tam tikrame algoritme. Kai sistema pasiima pirmą šabloną iš vartotojo sudarytų šablonų ir nustato šablono tipą, t.y. „*daiktavardinių frazių*“ šablonas ar „*veiksmožodinių frazių*“ šablonas, automatiškai pereina į atitinkamą veiksmą. Daiktavardinių frazių paieškos veiksmas padetalizuotas pateiktas (**Error! Reference source not found.**) veiklos diagramoje, veiksmožodinių frazių paieškos veiklos diagrama pateikta (29 Pav. Veiksmožodinių frazių atpažinimo veiklos diagrama) paveiksle. Jei sistemoje neįvyko klaidų beiškant frazių, tai sistema išsaugo visas rastas frazes ir pereina prie kito šablono. Kai sistema išanalizavo tekstą su visais frazių radimo šablonais, tai pereina prie rastų frazių apdorojimo veiksmo. Frazių papildomas apdorojimas susideda iš dviejų pagrindinių veiksmų:

- „*Atrinkti unikalias frazes*“ – šiame veiksmo programa automatiškai išrenka tik unikalias frazes.
- „*Atrinkti ilgiausias frazes*“ – šiame veiksmo yra išrenkamos ilgiausios frazės iš visų pasikartojančių ar įeinančių į tą frazę elementų, pvz.:
Frazių paieškos algoritmas rado tokias frazes: *namas, didelis, raudonas namas, didelis raudonas namas*. Algoritmas automatiškai išrenka tik ilgiausias tinkančias frazes, šiuo atveju liktų tik „*didelis raudonas namas*“ frazė.

Papildomas veiksmas po frazių apdorojimo vyksta „*Suskaičiuoti statistiką*“. Šiame veiksmo programa suskaičiuoja kiek iš viso buvo rasta šių frazių tekste.

3.1.1. Šablonų sudedamosios dalys

Šabloną, kaip ir minėta anksčiau privalo sudaryti būtent keturios dalys. Pagrindinė ir esminė šablonų sudedamoji dalis yra patys šablono elementai. Šabloną gali sudaryti 5 šablono elementai, kurie kiekvienas iš jų atskirai gali turėti savo dedamąsias dalis.

Žemiau pateiktame paveikslėlyje matyti pats duomenų modelis. Kaip atrodo, tiek visas ieškomos frazės struktūros šablonas, tiek jį sudarantys frazės struktūros šablono elementai:



Frazės struktūros šabloną gali sudaryti nuo vieno iki penkių frazės struktūros šablono elementų. Kiekvieną šablono elementą sudaro 10 dedamųjų:

- Kalbos dalis – žodžio kalbos dalis pvz.: *daiktavardis, būdvardis* ir t.t.
- Giminė – žodžio giminė pvz.: *moteriška, vyriška*.
- Linksnis – žodžio linksnis
- Laipsnis – žodžio laipsnis
- Asmuo – žodžio asmuo
- Laikas – žodžio laikas
- Pasikartojamumas – nusako ar šis šablono elementas gali kartotis, tol kol kitas elementas neatitinka jo struktūros.
- Privalomumas – nusako ar šis elementas yra privalomas šablone.
- Pažymimasis daiktavardis – nurodo ar šiame elemente esantis daiktavardis yra pažymimasis.
- Pažymimojo daiktavardžio dalis – nurodo kokie kiti šablono elementai yra priklausomi prie pažymimojo daiktavardžio.

Žemiau pateiktoje (9 Lentelė Kalbos dalių morfologinės savybės) lentelėje yra pateiktos lietuvių kalbos dalių, morfologinės kaitymo savybės. Kokioms kalbos dalims būdingas linksnis, giminė ir t.t.

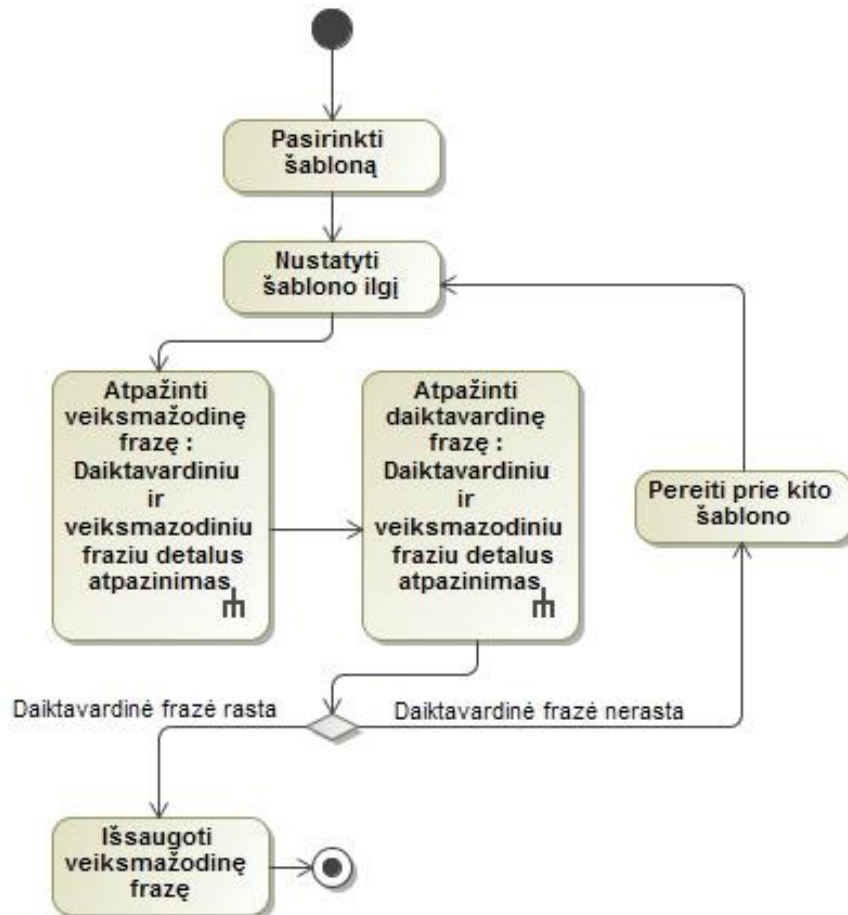
9 Lentelė Kalbos dalių morfologinės savybės

| Kalbos dalis | Morfologinės savybės |
|--------------|---|
| Daiktavardis | Skaičiais, linksniais, giminėmis |
| Būdvardis | Skaičiais, linksniais, giminėmis |
| Veiksmažodis | Asmenimis, skaičiais, laikais, nuosakomis |
| Prieveiksmis | ----- |
| Skaitvardis | Skaičiais, linksniais, giminėmis |
| Įvardis | Giminėmis, linksniais, skaičiais |
| Prielinksnis | ----- |
| Jungtukas | ----- |
| Jaustukas | ----- |
| Ištiktukas | ----- |
| Dalelytė | ----- |

3.1.2. Frazijų paieška naudojant šablonus

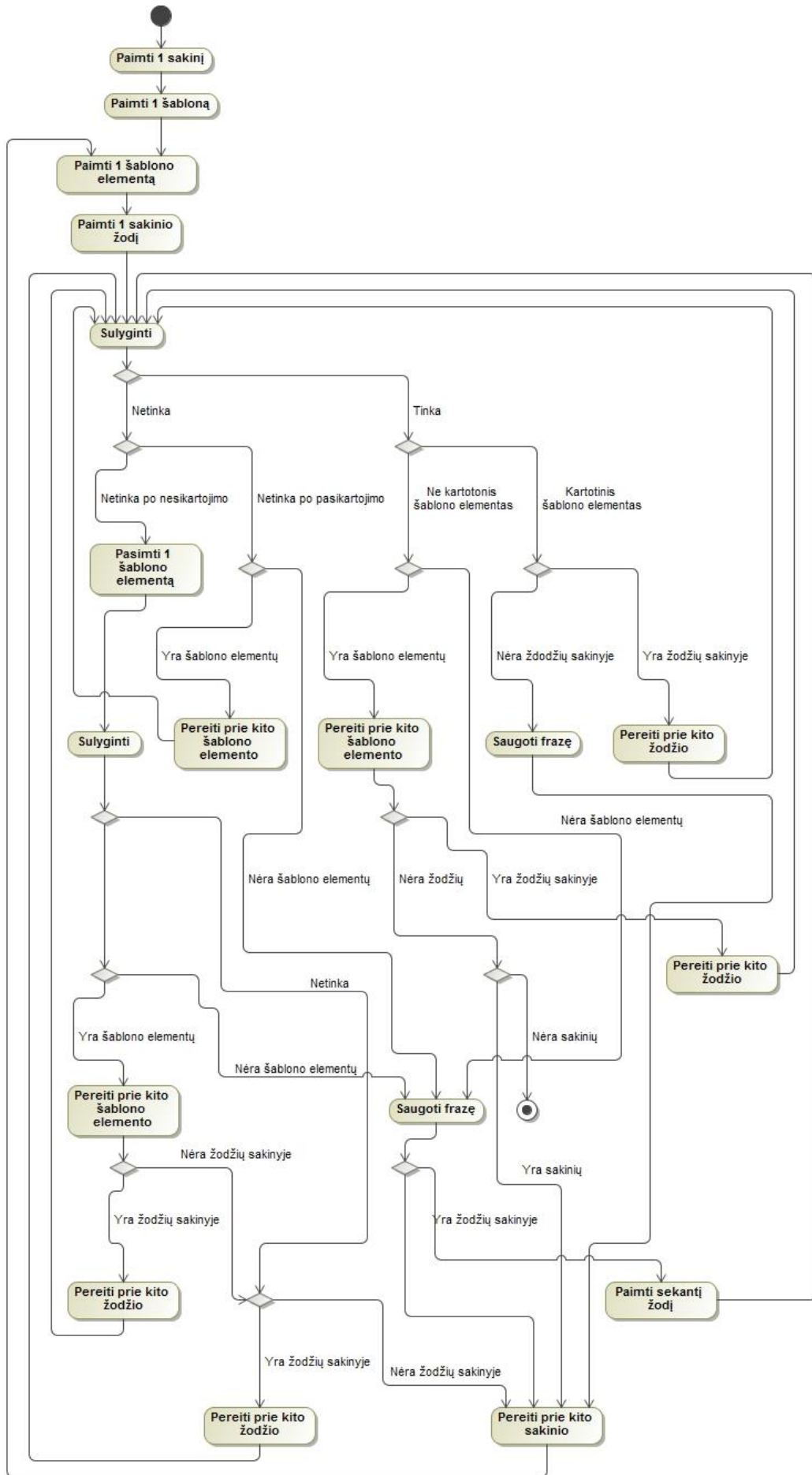
Kaip ir kiekvieni tekste esantys žodžiai taip ir šablono elementai turi tuos pačius elementus t.y. kalbos dalį, giminę, skaičių ir t.t. Teksto analizės metu frazės ieškomos nagrinėjant po vieną sakinį tikrinant ar esantys žodžiai su savo kalbos dalimis atitinka frazės struktūros šablono elementus. Jei visi esantys frazės struktūros elementai atitinka tekste esančių žodžių kalbos dalims ir eilės tvarka, tai šis žodžių junginys rastas tekste traktuojamas kaip frazė.

Veiklos diagramos „Frazijų radimo naudojant šablonus algoritmas“ kuris yra pateiktas (28 Pav. „Frazijų radimo naudojant šablonus algoritmas“ veiklos diagrama) paveiksle veiksmo „Atpažinti veiksmažodines frazes“ detalesnė veiklos diagrama pateikta žemiau (28 Pav. „Frazijų radimo naudojant šablonus algoritmas“ veiklos diagrama) paveiksle.



29 Pav. Veiksmažodinių frazių atpažinimo veiklos diagrama

Žemiau pateiktame (30 Pav. Detali daiktavardinių veiksmažodinių frazių paieškos algoritmo veiklos diagrama) paveiksle pavaizduota frazių paieškos veiklos diagrama. Ši veiklos diagrama detalizuoja pagrindinį frazių paieškos veiklos algoritmą, pagal kurį sistema ieško frazių naudodama šablonus. Šis algoritmas yra naudojamas tiek daiktavardiniams šablonams, tiek veiksmažodiniams.



30 Pav. Detali daiktvardinių veiksmažodinių frazių paieškos algoritmo veiklos diagrama

3.1.3. Algoritmo taikymo prielaidos ir situacijos

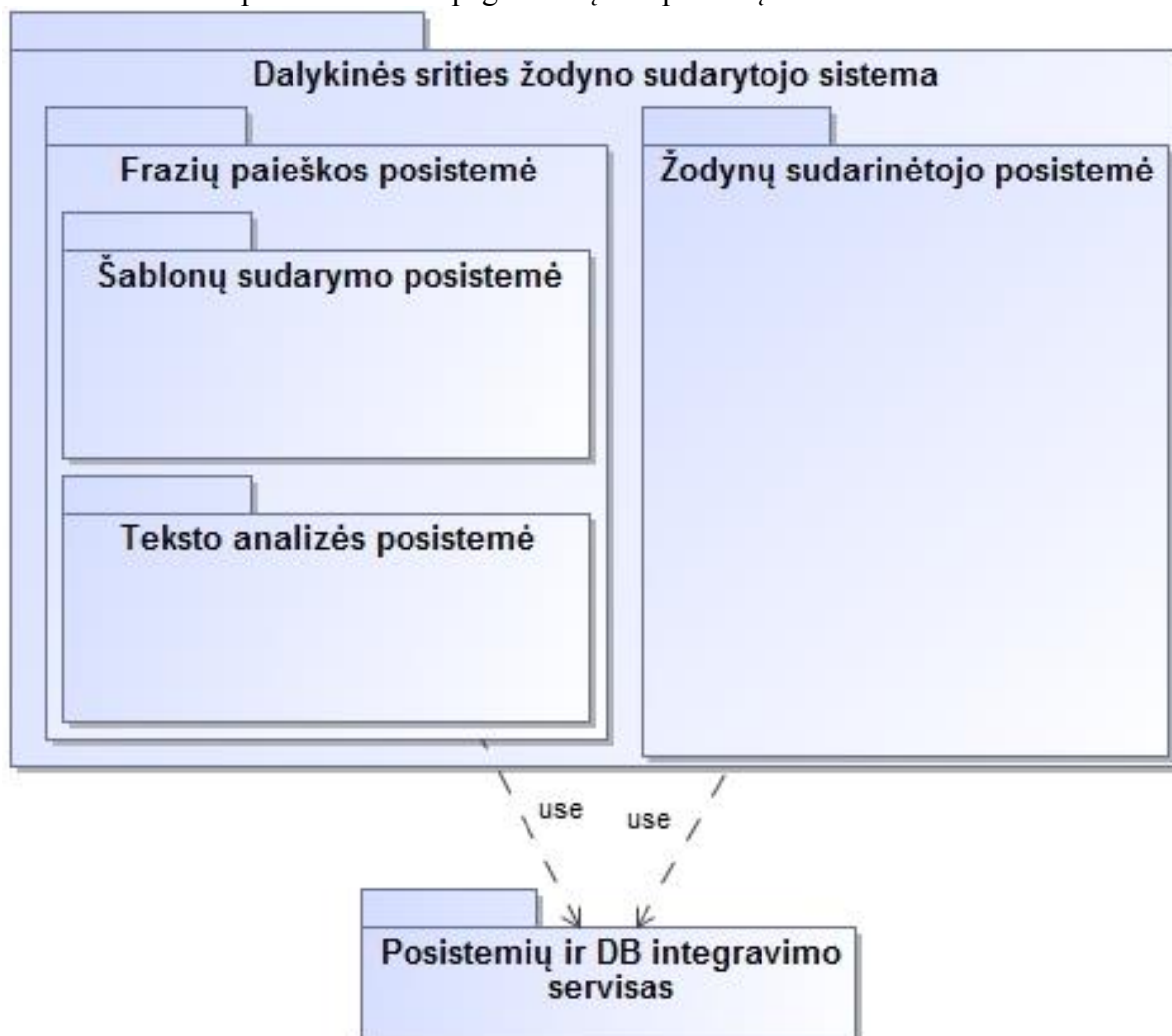
Šioje dalyje pateikiamos pagrindinės prielaidos, kurios gali lemti geresnių bei tinkamesnių frazių radimui. Taipogi pateikimas situacijos dėl kurių vartotojas gali būti apribotas.

- Vartotojas gali nenorėti vesti visų įmanomų frazių formų į sistemą, dėl to galimybė pasirinkti *pasikartojamumą* ir *privalomumą* frazės struktūros elementuose, leis jam laisviau ir lengviau sudaryti šablonus.
- Norint rasti kuo tinkamesnes frazes, būtina gerai ir kruopščiai sudaryti frazės struktūrų šablonus.
- Sudarinėjant šablonus ir juos skirstant į grupes bus galima rasti specifines, vartotojo norimas frazes, pagal jo užduotus frazių struktūrų šablonus.
- Algoritmas yra pritaikytas ieškoti penkių ilgio frazių, kadangi *praplėčiamumas* ir *pasikartojamumas* leidžia praplėsti paiešką nuo penkių iki dvidešimties žodžių frazėje.

4. AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE SUDARYTŲ ŠABLONŲ PAGRINDU REALIZACIJOS PROJEKTAS

4.1. Sistemos architektūra

Šioje dalyje pateikimas pagrindinis sistemos architektūrinis sprendimas. Žemiau esančiame (31 Pav. Dalykinės srities žodyno sudarytojo sistemos architektūrinis sprendimas) paveiksle pateikimas sistemos architektūrinis sprendimas tik iš pagrindinių komponentų.

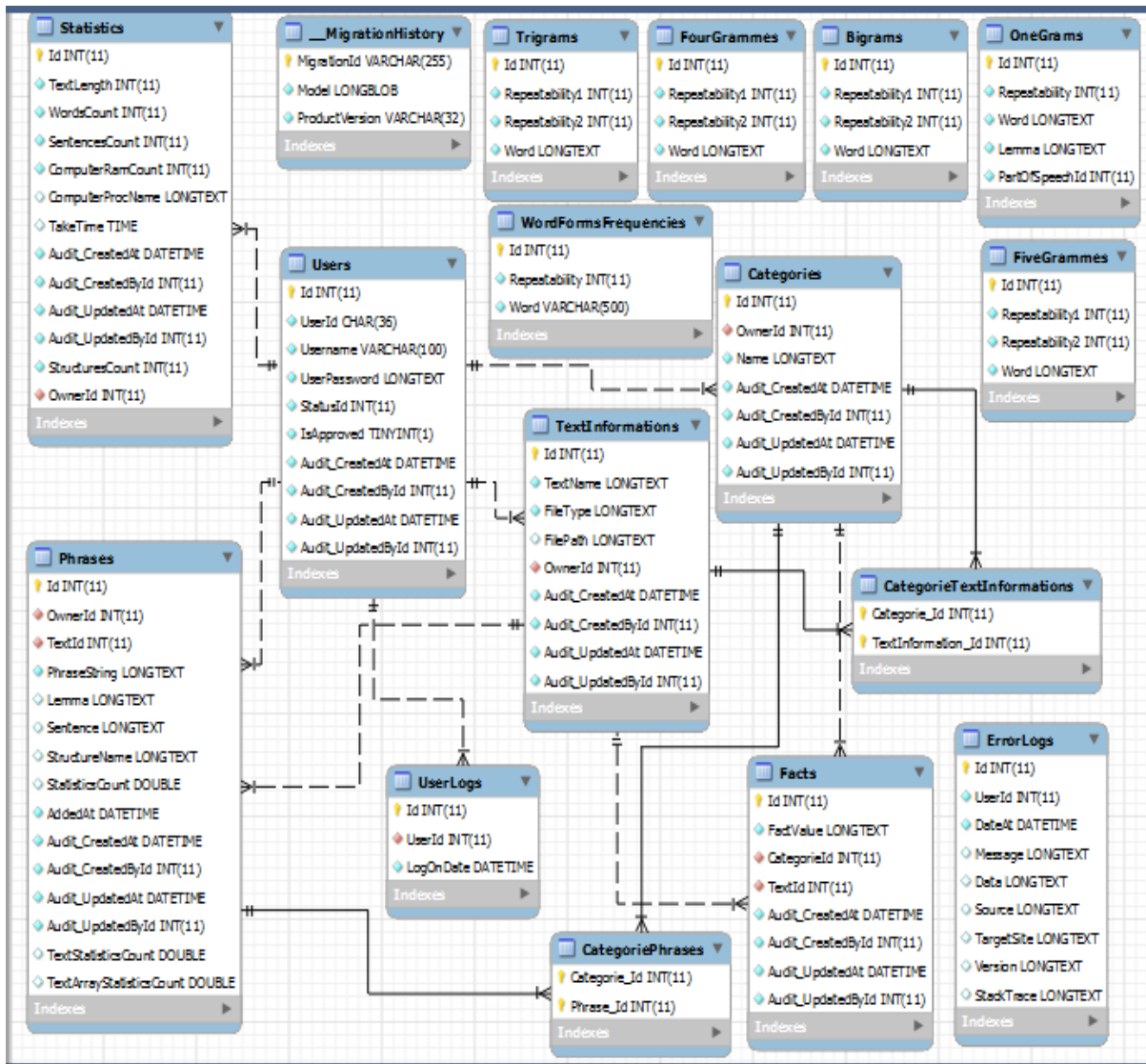


31 Pav. Dalykinės srities žodyno sudarytojo sistemos architektūrinis sprendimas

Dalykinės srities žodyno sudarytojo sistemą sudaro du pagrindiniai komponentai (posistemės), t.y. „Frazių paieškos posistemė“ ir „Žodynų sudarinėtojo posistemė“. Frazių paieškos posistemėje yra dar dvi papildomos posistemės kurios yra atsakingos už frazių šablonų sudarymą bei teksto analizės dalį. Šablonų sudarymo posistemės detalesnis vaizdas pateiktas (14 Pav. Šablonų sudarinėjimo posistemės panaudojimo atvejų diagrama) paveiksle, taipogi teksto analizės posistemės detalesnis vaizdas pateikiamas (20 Pav. Teksto analizės posistemės panaudojimo atvejų diagrama) paveiksle.

4.1.1. Duomenų bazės schema

Šioje duomenų bazėje yra saugoma visos vartotojų sudarytos frazės, vartotojo kategorijos, grupės. Duomenų bazėje taipogi saugoma visa informacija apie vartotojus, kaupiama jų statistika. Taipogi šioje duomenų bazėje yra sukaupta statistika iš kitų tekstų, kuri yra naudojama vartotojo rastų frazių statistikai skaičiuoti.

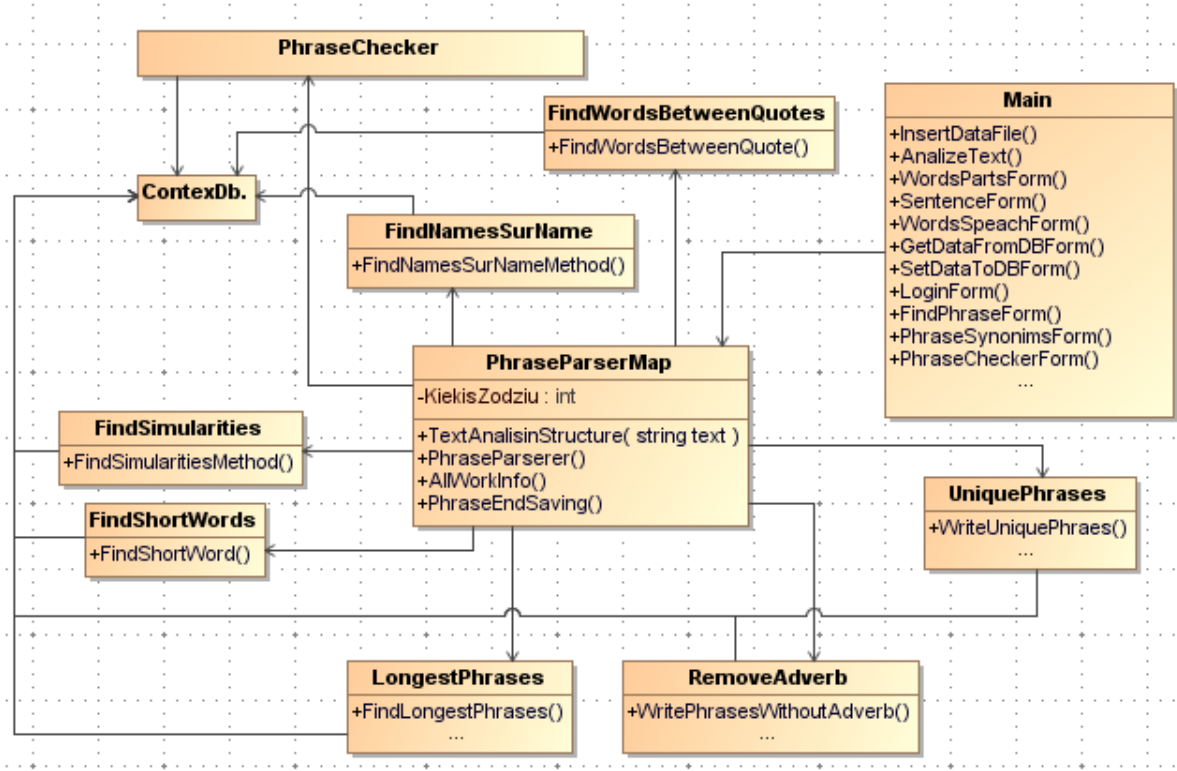


32 Pav. Pagrindinės duomenų bazės schema

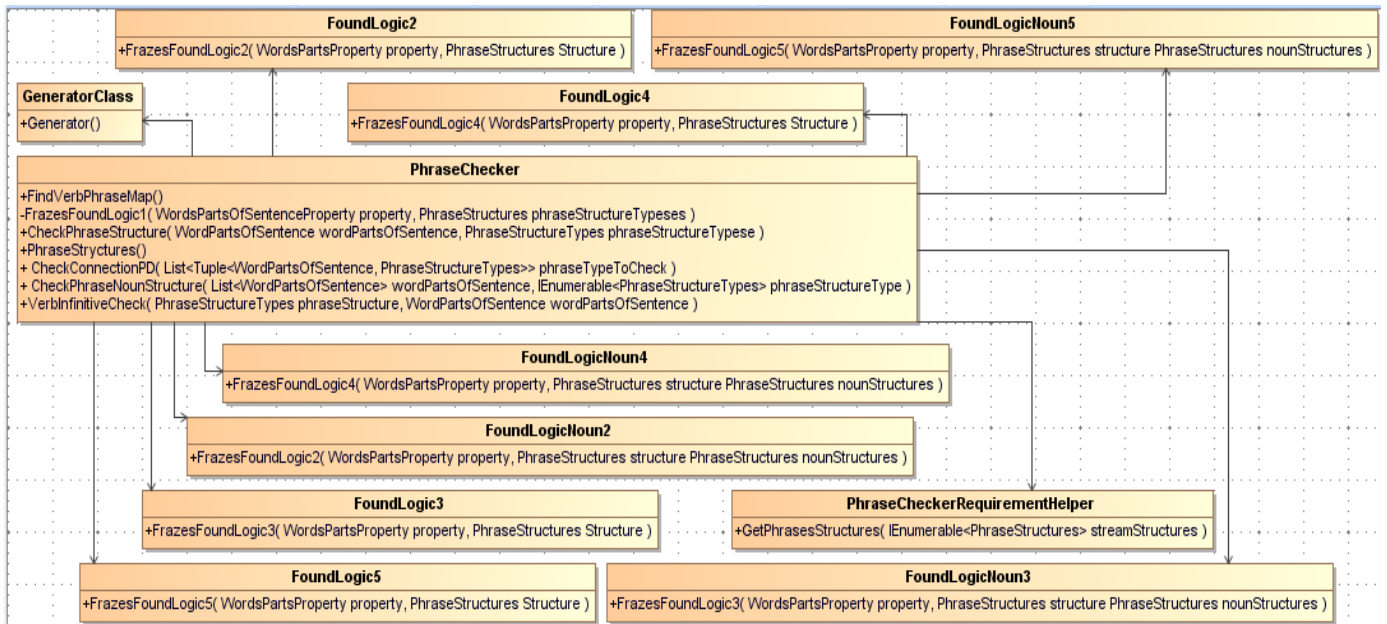
10 Lentelė Pagrindinės duomenų bazės schemos lentelių aprašai

| Lentelės pavadinimas | Aprašymas |
|----------------------|--|
| Users | Saugoma informacija apie vartotojus, kurie gali prisijungti prie serviso |
| PhraserTable | Saugomos visos frazės, kurias atsiuntė vartotojai |
| UsersLogs | Saugoma informacija apie vartotoju prisijungimus |
| Categories | Saugomos informacijos apie kategorijas |
| Statistics | Saugoma visa statistika |
| Facts | Saugomi vartotojo sudaryti faktai |
| ErrorLogs | Saugomi gauti klaidos pranešimai |

4.2. Frazių paieškos posistemės detalus projektas



33 Pav. Frazių radimo klasių daigrama



34 Pav. Frazių radimas naudojant šablonus klasių diagrama

Klasifikacija: Paketas.

Apibrėžimas: Pakete pateikiamos klasės kurios atsakingos už frazių radimą.

Atsakomybės: Paketas atsako už frazių radimą.

Struktūra: Komponentą sudaro klasės aprašytos pakete „Teksto analizės posistemė“ ir pateiktos diagramoje.

Sąveikavimas: Frazių radimo paketas naudojami klaidos paketu.

Resursai: Paketas naudoja lokalia DB.

Skaičiavimai: Skaičiavimai detalizuojami paketo klasių metodų aprašymuose.

Sąsaja/eksportas: Public metodai.

Klasė „FindWordsBetweenQuotes“

Apibrėžimas: Suranda frazes, kurios yra tarp kabučių.

Atsakomybės: Surasti frazes, kurios parašytos tarp kabučių.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *FindWordsBetweenQuotes* naudojami *contextDB* metodais, o *FindWordsBetweenQuotes* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindWordsBetweenQuote*

Klasės metodai

11 Lentelė FindWordsBetweenQuote metodas.

| Public FindWordsBetweenQuote | |
|------------------------------|---|
| <i>Atsakomybės</i> | Surasti frazes, kurios yra parašytos tarp kabučių. |
| <i>Skaičiavimai</i> | Nagrinėja tekstą sakiniiais ir ieško frazių tarp kabučių. |
| <i>Sąsaja/eksportas</i> | - |
| <i>Išimtys</i> | Klaidos atveju įrašas nėra įrašomas ir tęsiama paieška |

Klasė „ FindNamesSurName “

Apibrėžimas: Suranda vardus ir pavardes.

Atsakomybės: Surasti vardus ir pavardes iš teksto.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *FindNamesSurName* naudojami *contextDB* metodais, o *FindNamesSurName* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindNamesSurNameMethod*

Klasės metodai

12 Lentelė FindNamesSurNameMethod metodas.

| Public FindNamesSurNameMethod | |
|-------------------------------|---|
| <i>Atsakomybės</i> | Suranda vardus ir pavardes iš teksto. |
| <i>Skaičiavimai</i> | Nagrinėdami tekstą, pagal sakinius, ieškome žodžių kurie eina greta ir yra iš didžiųjų raidžių. |
| <i>Sąsaja/eksportas</i> | <i>DataBaseContext</i> |
| <i>Išimtys</i> | Jai paskutinis žodis arba sakinyis iš 1 žodžio. |

Klasė „FindShortWords“

Apibrėžimas: Suranda sutrumpinimus.

Atsakomybės: Surasti sutrumpinimus iš teksto (pvz. KTU).

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *FindShortWords* naudojami *contextDB* metodais, o *FindShortWords* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindShortWord*

Klasės metodai

13 Lentelė FindShortWord metodas.

| <i>Public FindShortWord</i> | |
|-----------------------------|---|
| <i>Atsakomybės</i> | Surasti sutrumpinimu, tokius kaip KTU |
| <i>Skaičiavimai</i> | Ieškomi žodžiai kurie yra sudaryti vien iš didžiųjų raidžių. |
| <i>Sąsaja/eksportas</i> | Rasti sutrumpinimai įrašomi į duomenų bazę naudojant <i>WritePhrase</i> metodą. |
| <i>Išimtys</i> | - |

Klasė „FindSimilarities“

Apibrėžimas: Suranda sinoniminius žodžius.

Atsakomybės: Suranda sinoniminius žodžius.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *FindSimilarities* naudojami *contextDB* metodais, o *FindSimilarities* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindSimilaritiesMethod*

Klasės metodai

14 Lentelė FindSimilaritiesMethod metodas.

| <i>Public FindSimilaritiesMethod</i> | |
|--------------------------------------|---|
| <i>Atsakomybės</i> | Suranda sinoniminius žodžius tarp sutrumpintų frazių ir teksto. |
| <i>Skaičiavimai</i> | Ieško atitikmenų tekste su sutrumpintomis frazėmis. |
| <i>Sąsaja/eksportas</i> | Duomenų įrašymui į lokalią duomenų bazę <i>contextDB</i> |
| <i>Išimtys</i> | - |

Klasė „LongestPhrases“

Apibrėžimas: Suranda ilgiausias frazes.

Atsakomybės: Palikti ilgiausias frazes ir pašalinti trumpesnes, kurios įeina į ilgesne frazes.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *LongestPhrases* naudojami *contextDB* metodais, o *LongestPhrases* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindLongestPhrases*

Klasės metodai

15 Lentelė FindLongestPhrases metodas.

| Public FindLongestPhrases | |
|---------------------------|---|
| <i>Atsakomybės</i> | Ilgiausių frazių radimas |
| <i>Skaičiavimai</i> | Randamos ilgiausios rastos frazės sakinyje. |
| <i>Sąsaja/eksportas</i> | Paimamos rastos frazės iš lokalių duomenų bazės. Įrašomos tik ilgiausios rastos frazės į <i>contextDB</i> . |
| <i>Išimtys</i> | - |

Klasė „RemoveAdverb“

Apibrėžimas: Pašalina frazių priekyje einančius jungtukus, prielinksnius.

Atsakomybės: Pašalina frazių priekyje einančius jungtukus, prielinksnius.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *RemoveAdverb* naudojami *contextDB* metodais, o *RemoveAdverb* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *WritePhrasesWithoutAdverb*

Klasės metodai

16 Lentelė WritePhrasesWithoutAdverb metodas.

| Public WritePhrasesWithoutAdverb | |
|----------------------------------|---|
| <i>Atsakomybės</i> | Pašalinti frazių priekyje einančius jungtukus, prieveiksmius. |
| <i>Skaičiavimai</i> | Iš visų rastų frazių pašalinami pirmi žodžiai kurie yra arba jungtukai, arba prieveiksmiai. |
| <i>Sąsaja/eksportas</i> | - |
| <i>Išimtys</i> | - |

Klasė „ UniquePhrases “

Apibrėžimas: Suranda iš visų rastų frazių tik unikalias frazes.

Atsakomybės: Suranda iš visų rastų frazių tik unikalias.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *UniquePhrases* naudojami *contexDB* metodais, o *UniquePhrases* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *WriteUniquePhraes*

Klasės metodai

17 Lentelė *WriteUniquePhraes* metodas.

| Public <i>WriteUniquePhraes</i> | |
|---------------------------------|--|
| <i>Atsakomybės</i> | Surasti iš visų surastų frazių tik unikalias frazes |
| <i>Skaičiavimai</i> | Gaunamos visos frazes, einama po vieną fraze ir ieškoma ar tokia jau yra sąrašė, jai yra nerašome jai nėra įrašoma i nauja unikaliu frazių sąrašą. |
| <i>Sąsaja/eksportas</i> | - |
| <i>Išimtys</i> | - |

Klasė „ PhraseChecker“

Apibrėžimas: Ieško daiktavardinių ir veiksmažodinių frazių pagal šablonus.

Atsakomybės: Ieško daiktavardinių ir veiksmažodinių frazių pagal šablonus

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *PhraseChecker* naudojami *contexDB*, *FoundLogic2*, *FoundLogic3*, *FoundLogic4*, *FoundLogic5*, *GeneratorClass*, *FoundLogicNoun2*, *FoundLogicNoun3*, *FoundLogicNoun4*, *FoundLogicNoun5*, *PhraseCheckerRequirementHelper* metodais, o *PhraseChecker* metodais naudojami *PhraseParserMap*.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FindVerbPhraseMap*, *FrazesFoundLogic1*, *CheckPhraseStructure*, *CheckPhraseStructure*

Klasės metodai

18 Lentelė *FindVerbPhraseMap* metodas.

| Public <i>FindVerbPhraseMap</i> | |
|---------------------------------|--|
| <i>Atsakomybės</i> | Ieško daiktavardiniu ir veiksmažodinių frazių. |
| <i>Skaičiavimai</i> | Gaunami pateikti šablonai taipogi sakiniai su tam sakiniui priklausančiais žodžiais ir jų kalbos dalimis. Tikriname sakinius tekstą naudojant gautus šablonus (daiktavardinį ar veiksmažodinį šabloną) |
| <i>Sąsaja/eksportas</i> | - |
| <i>Išimtys</i> | Nėra sudaryta šablonų. |

19 Lentelė FrazesFoundLogic1(WordsPartsOfSentenceProperty property, PhraseStructures phraseStructureTypeses) metodas.

| | |
|---|--|
| <i>Public</i> FrazesFoundLogic1(WordsPartsOfSentenceProperty property, PhraseStructures phraseStructureTypeses) | |
| <i>Atsakomybės</i> | Suranda daiktavardines frazes iš vieno šablono elemento |
| <i>Skaičiavimai</i> | Ieškoma sakinyje pagal einančių žodžių kalbos dalis aitkmenų šablone |
| <i>Sąsaja/eksportas</i> | Įrašomos rastos frazės į sąrašą. |
| <i>Išimtys</i> | Jai nėra sudarytų šablonų, skaičiavimai nevyksta. |

20 Lentelė CheckPhraseStructure(WordPartsOfSentence wordPartsOfSentence, PhraseStructureTypes phraseStructureTypeese) metodas.

| | |
|--|---|
| <i>Public</i> CheckPhraseStructure(WordPartsOfSentence wordPartsOfSentence, PhraseStructureTypes phraseStructureTypeese) | |
| <i>Atsakomybės</i> | Patikrina ar žodis atitinka frazes šablono elementą. |
| <i>Skaičiavimai</i> | Patikrina gauto žodžio struktūra su frazių šablono struktūra. |
| <i>Sąsaja/eksportas</i> | Grąžina <i>true</i> – jei elementas tinkamas, arba <i>false</i> – jai elementas netinkamas. |
| <i>Išimtys</i> | - |

21 Lentelė PhraseStructures metodas.

| | |
|--------------------------------|--|
| <i>Public</i> PhraseStructures | |
| <i>Atsakomybės</i> | Gauna visas struktūras, kurias sudarė vartotojas, jai vartotojas nesudarė fazių struktūrų, tai nauji šablonai sugeneruojami automatiškai sistemos. |
| <i>Skaičiavimai</i> | - |
| <i>Sąsaja/eksportas</i> | Grąžina <i>PhraseStructures</i> sąrašą. |
| <i>Išimtys</i> | - |

22 Lentelė CheckConnectionPD(List<Tuple<WordPartsOfSentence, PhraseStructureTypes>> phraseTypeToCheck) metodas.

| | |
|---|---|
| <i>Public</i> CheckConnectionPD(List<Tuple<WordPartsOfSentence, PhraseStructureTypes>> phraseTypeToCheck) | |
| <i>Atsakomybės</i> | Patikrina ar daiktavardis yra pažymimasis. |
| <i>Skaičiavimai</i> | Patikrina atėjusioje struktūroje ar yra pažymimasis daiktavardis bei prie PD yra elementai kurie turi jį atitikti. Tikrinama ar žodžių masyvas atitinka pažymimojo daiktavardžio ir jį atitinkančių elementų kombinaciją. Jai <i>true</i> – rastas pažymimasis daiktavardis, jei <i>false</i> – ieškoma toliau pagal sekančius žodžius, tikrinant pagal PD. |
| <i>Sąsaja/eksportas</i> | Grąžina <i>true</i> arba <i>false</i> |
| <i>Išimtys</i> | Jai šablone nėra pažymimojo daiktavardžio, šis metodas nevykdomas. |

23 Lentelė VerbInfinitiveCheck(PhraseStructureTypes phraseStructure, WordPartsOfSentence wordPartsOfSentence)metodas.

| | |
|--|--|
| <i>Public</i> VerbInfinitiveCheck(PhraseStructureTypes phraseStructure, WordPartsOfSentence wordPartsOfSentence) | |
| <i>Atsakomybės</i> | Patikrina ar veiksmažodis yra bendratis. |
| <i>Skaičiavimai</i> | - |
| <i>Sąsaja/eksportas</i> | Grąžina <i>true</i> arba <i>false</i> , <i>DataBaseContext</i> |
| <i>Išimtys</i> | Jai ne bendratis, tačiau kalbos dalis yra <i>Veiksmažodis</i> , tai iš duomenų bazės randama bendratis atitinkantį veiksmažodžio žodį. |

Klasė „ FoundLogic2“

Apibrėžimas: Ieško daiktavardinių frazių, kurie sudaryti iš dviejų elementų.

Atsakomybės: Surasti daiktavardines frazes sudarytas iš dviejų elementų.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *PhraseChecker* naudojami *FoundLogic2* metodais.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FrazesFoundLogic2*

Klasės metodai

24 Lentelė FrazesFoundLogic2metodas.

| | |
|--|--|
| <i>Public</i> FrazesFoundLogic2(WordsPartsProperty property, PhraseStructures Structure) | |
| <i>Atsakomybės</i> | Ieško daiktavardinių frazių kurios atitinka dviejų elementų sudarytus šablonus. |
| <i>Skaičiavimai</i> | Gauna sakinio žodžius su jų kalbos dalimis, vėliau tikrinama ar šios kalbos dalys atitinka pateiktą daiktavardinės frazės šablono struktūrą. |
| <i>Sąsaja/eksportas</i> | Grąžina rastų frazių sąrašą. |
| <i>Išimtys</i> | - |

Klasės *FrazesFoundLogic3*, *FrazesFoundLogic4*, *FrazesFoundLogic5* logika bei metodai atitinka *FrazesFoundLogic2* logika, tik ieškoma pagal atitinkamai ilgesnius šablonus.

Klasė „GeneratorClass“

Apibrėžimas: Sugeneruoja standartinius frazių radimo šablonus.

Atsakomybės: Sugeneruoja standartinius frazių radimo šablonus.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *PhraseChecker* naudojami *FoundLogic2* metodais.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *Generator*

Klasės metodai

25 Lentelė *Generator* metodas.

| <i>Public Generator</i> | |
|-------------------------|---|
| <i>Atsakomybės</i> | Sugeneruoti standartinių šablonų sąrašą. |
| <i>Skaičiavimai</i> | Pagal tam tikras taisykles sugeneruoja frazių šablonus. |
| <i>Sąsaja/eksportas</i> | Gražina frazių šablonus. |
| <i>Išimtys</i> | - |

Klasė „PhraseCheckerRequirementHelper“

Apibrėžimas: Generuoti naujas frazių struktūras pagal vartotojo šablonus.

Atsakomybės: Generuoti naujas frazių struktūras pagal vartotojo šablonus.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *PhraseChecker* naudojami *PhraseCheckerRequirementHelper* metodais.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *GetPhrasesStructures*

Klasės metodai

26 Lentelė *GetPhrasesStructures* metodas.

| <i>Public GetPhrasesStructures(IEnumerable<PhraseStructures> streamStructures)</i> | |
|--|--|
| <i>Atsakomybės</i> | Sugeneruoti naujus frazių šablonų struktūras. |
| <i>Skaičiavimai</i> | Iš pravalomų ir neprivalomų kalbos dalių elementų sudaromos naujos visos įmanomos frazių struktūros. |
| <i>Sąsaja/eksportas</i> | Gražinama visos galimos frazių struktūros. |
| <i>Išimtys</i> | - |

Klasė „ FoundLogicNoun2 “

Apibrėžimas: Ieško visų veiksmažodinių frazių sudarytų iš dviejų elementų.

Atsakomybės: Surasti visas veiksmažodines frazes sudarytas iš dviejų elementų.

Struktūra: Klasės struktūra pateikta klasių diagramoje.

Sąveikavimas: *PhraseChecker* naudojami *FoundLogicNoun2* metodais.

Resursai: Klasė nenaudoja jokių papildomų resursų.

Skaičiavimai: Skaičiavimai pateikti klasės metodų aprašymuose

Sąsaja/eksportas: Metodai: *FrazesFoundLogic2*

Klasės metodai

27 Lentelė FoundLogicNoun2 metodas.

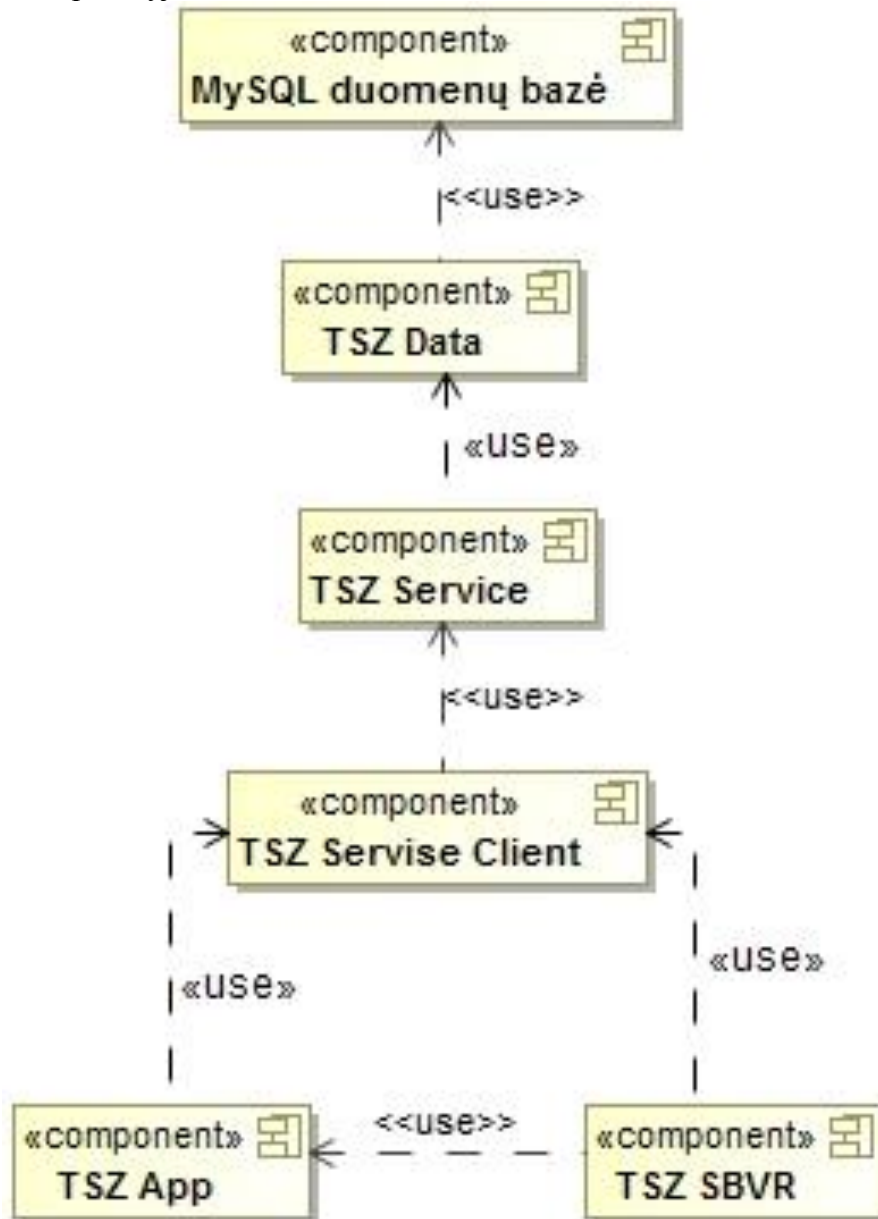
| <i>Public FoundLogicNoun2(WordsPartsProperty property, PhraseStructures structure PhraseStructures nounStructures)</i> | |
|--|--|
| <i>Atsakomybės</i> | Ieško veiksmažodinių frazių kurios atitinka dviejų elementų sudarytus šablonus. |
| <i>Skaičiavimai</i> | Gauna sakinio žodžius su jų kalbos dalimis, vėliau tikrinama ar šios kalbos dalys atitinka pateiktą veiksmažodinės frazės šablono struktūrą. |
| <i>Sąsaja/eksportas</i> | Rastų frazių sąrašas. |
| <i>Išimtys</i> | - |

Klasės *FoundLogicNoun3*, *FoundLogicNoun 4*, *FoundLogicNoun 5* logika bei metodai atitinka *FoundLogicNoun2* logika, tik ieškoma pagal atitinkamai ilgesnius šablonus.

4.3. Realizacijos modelis

4.3.1. Programinių komponentų architektūra

Žemiau pateiktame (35 Pav. Programų komponentų architektūra) paveiksle pateikta programos komponentų sąveika tarpusavyje.



35 Pav. Programų komponentų architektūra

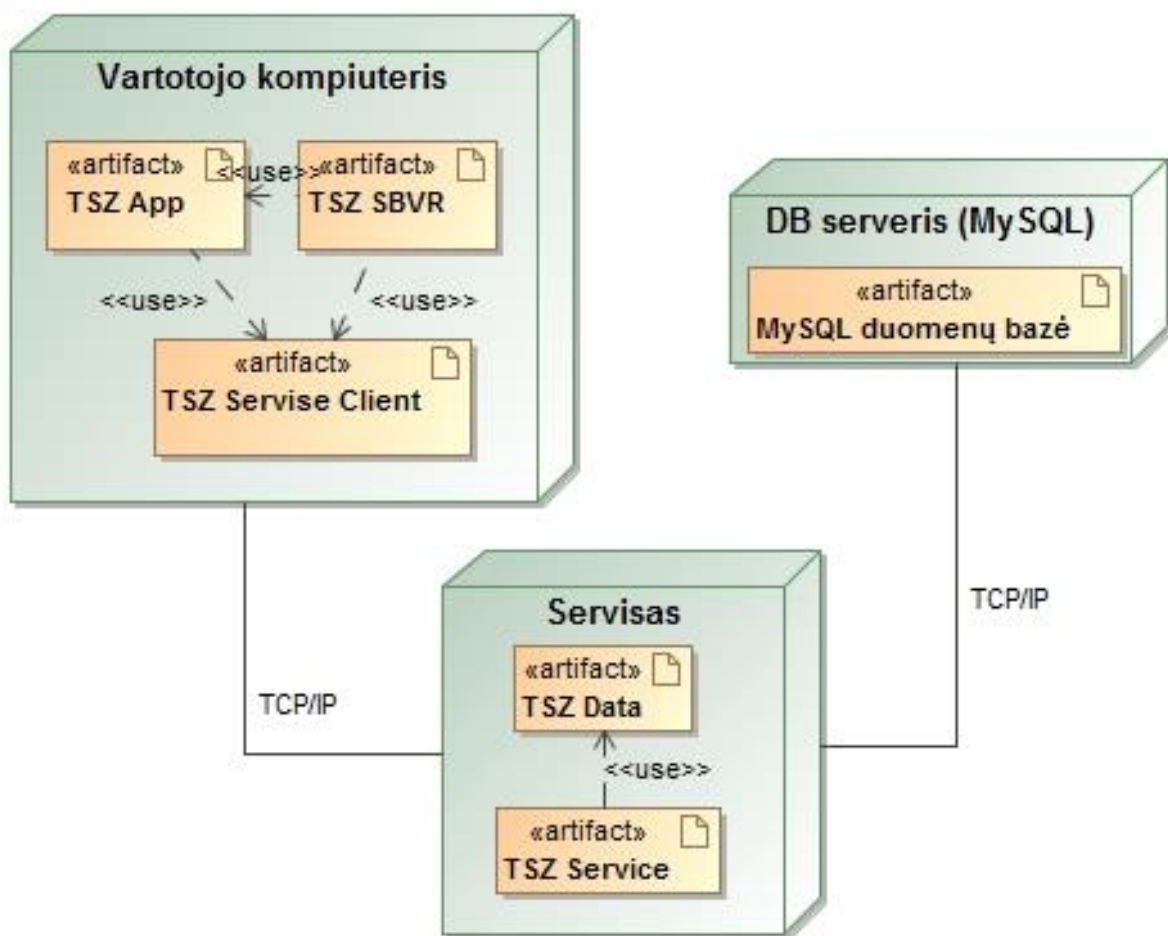
Sistemą sudaro šeši komponentai. Kiekvieno komponento trypas aprašymas pateiktas (28 Lentelė Komponentų paaiškinimų lentelė) lentelėje.

28 Lentelė Komponentų paaiškinimų lentelė

| Komponento pavadinimas | Paiškinimas |
|------------------------|--|
| TSZ App | Komponentas kuris atlieka pagrindinį programos veikimą: kuria šablonus, ieško frazių ir t.t. |
| TSZ SBVR | Komponentas atsakingas už SBVR žodyno sudarymą, |
| TSZ Service Client | Komponentas atsakingas už bendravimą su servisu bei duomenų siuntimą per TCP/IP. |
| TSZ Service | Komponentas atliekantis statistikos analizę, duomenų išsaugojimą. |
| TSZ Data | Komponentas kuris atsakingas už duomenų išsaugojimą į duomenų bazę. |
| MySQL duomenų bazė | Komponentas skirtas duomenims saugoti (duomenų bazė). |

4.3.2. Diegimo modelis

Žemiau pateiktame (36 Pav. Sistemos diegimo modelis) paveiksle vaizduojamas pagrindinis komponentų išsidėstymas diegimo būsenoje. Komponentai: *TSZ App*, *TSZ SBVR*, *TSZ Service Client* yra diegiami pas vartotoją į kompiuterį. Komponentai: *TSZ Data*, *TSZ Service* yra diegiami atskirai į nutolusį servisą. Kaip jau minėta komponentų architektūros dalyje, *TSZ Service Client* ir *TSZ Service* bendrauja TCP/IP protokolu, per internetą. Komponentas *MySQL duomenų bazė* yra duomenų bazė kurioje saugomi duomenys iš vartotojo kompiuterio, kurie pateikiami per servisą. Servisas su duomenų baze bendrauja internetu TCP/IP protokolu.



36 Pav. Sistemos diegimo modelis

5. AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE SUDARYTŲ ŠABLONŲ PAGRINDU REALIZACIJA

5.1. Realizacijos ir veikimo aprašymas

5.1.1. Darbo pradžia

Prieš pradėdant darbą sistemoje - Teminės srities žodyno sudarinėtojuje (toliau TSZ), reikia instaliuoti papildomus priedus, jai jų nėra jau instaliuota kompiuteryje:

- .NET 4 – 4,5
- Windows PowerPack
- SQL ServerPack

Pradėjus instaliuotis, mūsų sistemos instaliacijos vedlys pats pamatys, kokių priedų trūksta sistemoje, kad TSZ veiktų tinkamai ir korektiškai sistemoje.

5.1.1.1. Instaliavimas

Nuėjus į vietą, kur įkėlėte programos instaliavimo failus reikia paleisti failą „setup“.

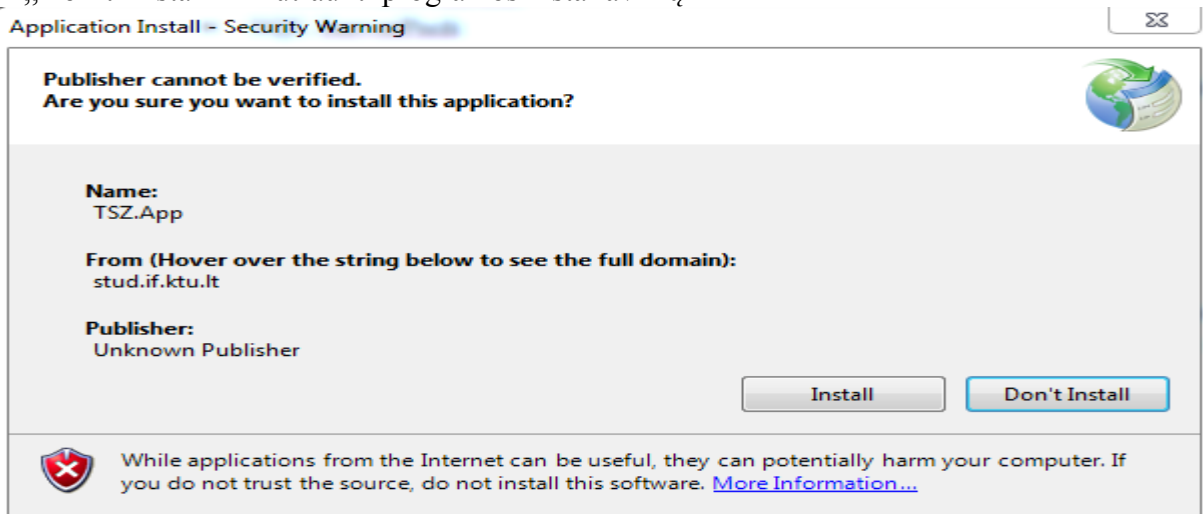
| | | | |
|---------|--------------------|----------------------|--------|
| setup | 12/15/2013 2:35 PM | Application | 847 KB |
| TSZ.App | 12/15/2013 2:58 PM | ClickOnce Applica... | 6 KB |

37 Pav. Instaliavimo failai

Atsiradusiame naujame lange (62 pav.) galimi du variantai : „Install“ ir „Don't Install“.

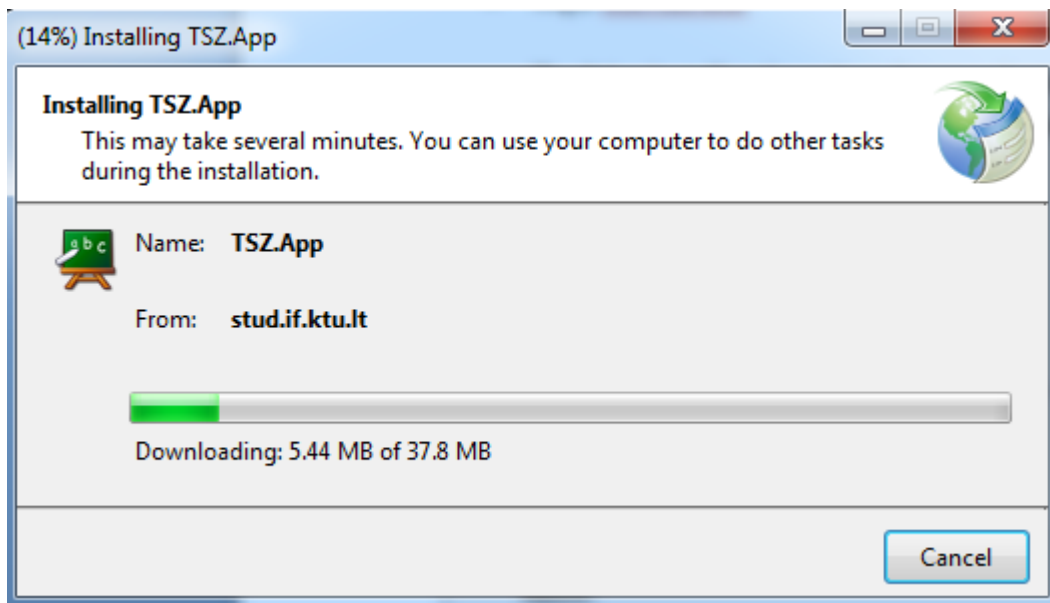
„Install“ – skirtas toliau tęsti programos instaliavimą

„Don't Install“ - nutraukti programos instaliavimą



38 Pav. Instaliavimo langas

Pasirinkus „Install“ toliau atsiranda instaliavimo langas (63 pav.), kuriame matome instaliavimo progresą.

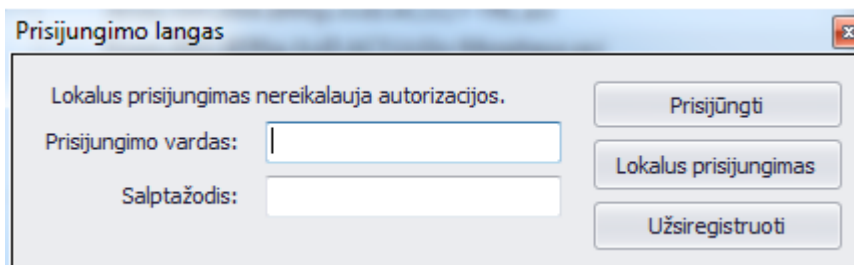


39 Pav. Instaliavimo progreso langas

Baigus instaliaciją, programa pirmą kartą pasileidžia automatiškai.

5.1.1.2. Prisijungimas

Kiekvieną kartą paleidžiant programa atsiranda prisijungimo langas (61 pav.). Prisijungimo lange galima pasirinkti, kuriame režime norima dirbti. Galima suvesti prisijungimo vardą ir slaptažodį ir paspaudus mygtuką „Prisijungti“ įgalinti funkcijas skirtas darbui su servisu. Lokalus prisijungimas turi visas funkcijas, išskyrus bendravimą su servisu. Užsiregistruoti skirta, naujam vartotojui prisiregistruoti prie serviso.

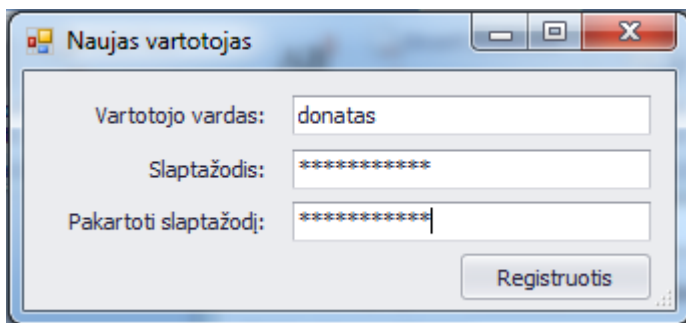


40 Pav. Prisijungimo langas

5.1.1.3. Registracija

Registracija skirta vartotojams, kurie norės naudotis programos servisu. Prisijungimo lange (65 pav.) reikia paspausti mygtuką „Užsiregistruoti“. Atsiradusiame naujame lange reikia užpildyti tokius laukus:

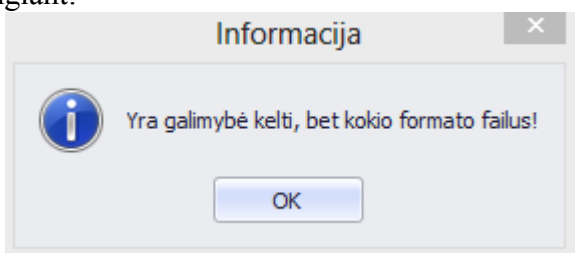
- Vartotojo vardas – vartotojo vardas, su kuriuos bus prisijungiama prie programos serviso.
- Slaptažodis – vartotojo slaptažodis, su kuriuo bus jungiamasi prie programos serviso
- Pakartoti slaptažodį – patikrinti ar vartotojas gerai suvedė slaptažodį



41 Pav. Registracijos langas

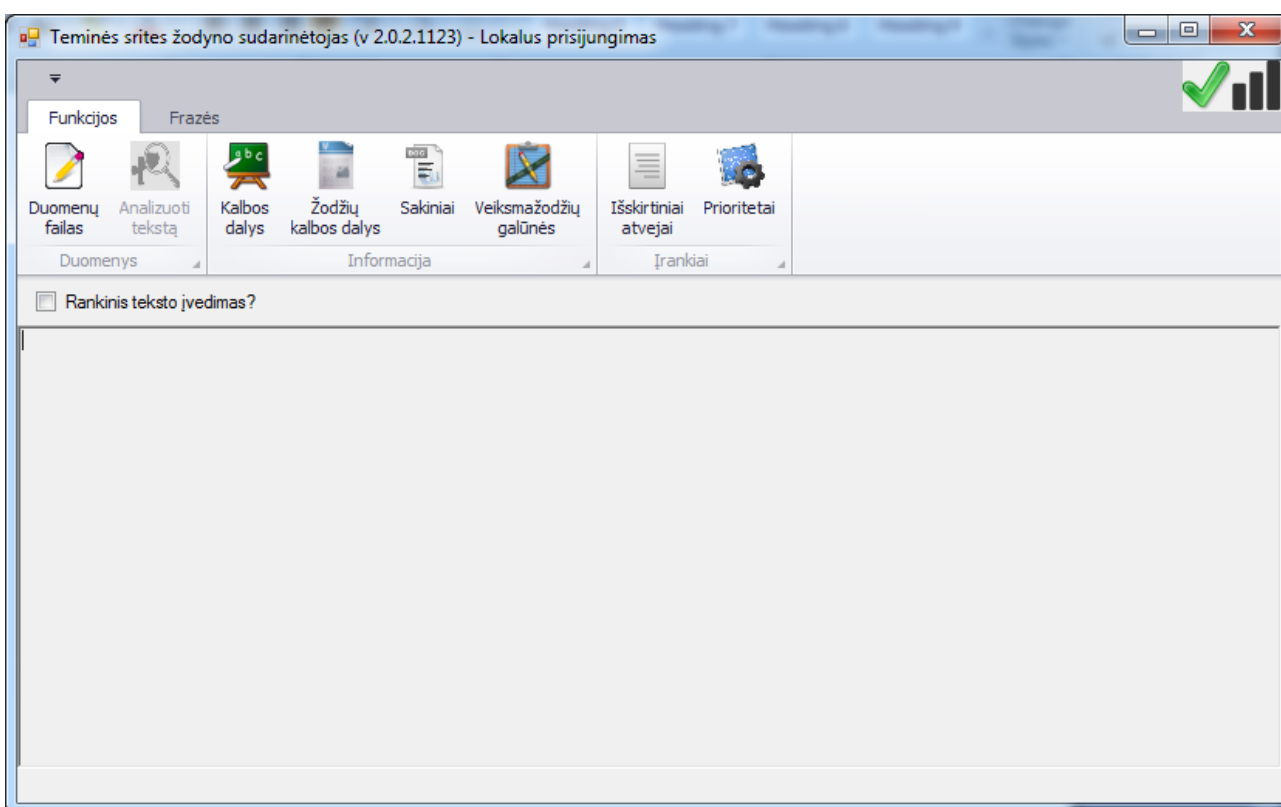
5.1.1.4. Sistemos informacija

Apie naujas sistemos funkcijas, ar naujus priedus, vartotojas yra supažindinamas, prieš sistemai įsijungiant:



42 Pav. Sistemos informacinis langas

5.1.2. Sistemos langas



43 Pav. Pagrindinis programos langas

Šis langas yra pagrindinis vartotojo funkcionalumo langas, kuriame jis galės pilnai valdyti TSZ sistemą.

5.1.2.1. Informaciniai ženklai

TSZ sistema turi du pagrindinius sistemos būklės nusakymo elementus:

- Serviso jungiamumo:

Ryšys su servisu yra:



Ryšio su servisu nėra:



- Interneto jungiamumo:

Internetas yra:



Interneto nėra:



Nesant ryšio su servisu nėra galimybės rasti frazes išsaugoti sistemos pagrindiniame serveryje. Nesant ryšio su internetu nėra galimybės kelti bet kokio tipo tekstą, kadangi nėra galimybės susisiekti su pagrindiniu VDU servisu, kuris galėtų morfologiškai išanotuoti įkeltą tekstą.

5.1.2.2. Sistemos funkciniai laukai

Pagrindinis sistemos langas yra suskirstytas į tris pagrindines dalis:

- Funkcijos – čia pateiktos pagrindinės naudojamos funkcijos, kurios yra dar papildomai suskirstytos į tris pogrupius: *Duomenys*, *Informacija*, *Įrankiai*.
- Frazės – čia pateiktas funkcijų sąrašas kurių pagalba galima valdyti frazių paiešką, tvarkyti rasti frazes, ar tiesiog peržiūrėti rastų frazių statistika (informacija: norint peržiūrėti frazių statistika būtina burėti ryšį su servisu). Ši dalis suskirstyta taipogi į tris pogrupius: *Frazės*, *Frazių statistika*, *Sinchronizavimas*.
- Teksto įvedimo langas.

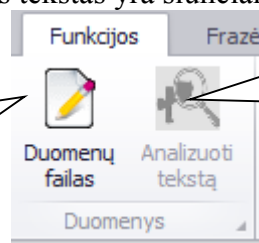
5.1.2.3. Duomenų pateikimas sistemai

Sistemai pateikti duomenis galima pateikti keliais būdais.

Pirmas būdas yra tiesiog sistemai pateikti suanotuotą duomenų failą. Šitokiam duomenų failo tipui nereikia interneto prisijungimo.

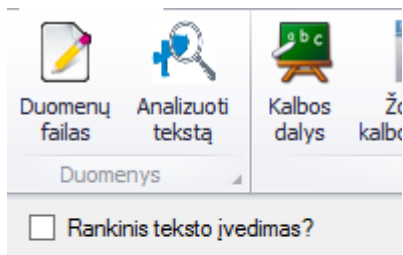
Sekantis būdas galima įkelti kitų formatų tekstus: *.pdf*, *.doc*, *.docx*, *.txt*, tačiau šiems formatams įkėlus

Duomenų failo įkėlimas, įkėlus duomenų failą automatiškai įsigalina „Teksto analizavimo“ funkcija



„Analizuoti tekstą“ funkcija kurios pagalba analizuojamas tekstas.

Įkėlus duomenų failą, galima pasirinkti iš karto analizuoti tekstą, arba galima šiek tiek vėliau, nuspaudus „Analizuoti tekstą“ mygtuką.

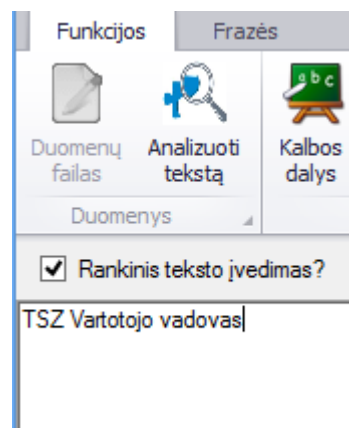


Sėkmingai įkėlus tekstą sistema įsigalina „Analizuoti tekstą“ mygtukas, taipogi analizuojamas tekstas yra pateikimas vartotojui, peržiūrėti žemiau esančiame lange.

Trečias būdas tiesiog įgalinus „Rankinio teksto“ vedimo funkciją, galima tekstą suvesti ranka. Tačiau čia kaip ir antru būdu privalomas internetas teksto suanotavimui VDU servise.

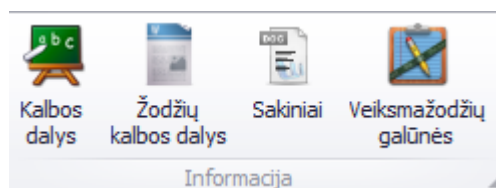
Keturi metai darbo, triskart keitėsis pavadi
Šis atviro pasaulio žaidimas pasakoja apie
organizuoto nusikalstamumo grupuotes, ži
bosus. Darbą nuo pat pradžių apsunkin
Darbą nuo pat pradžių apsunkina kelios a
antra, Wei Shen yra įsitikinęs, jog viena H
agentas netaptų savo vaikystės bičiuliu „
sąskaitų suvedinėjimo, žaidimo eigoje tai
norima išanalizuoti. Šio teksto analizavimui atlikti galima
nuspaudus „Analizuoti funkcija“.

Įgalinus rankinio teksto vedimą, galima parašyti tekstą, kurį



5.1.2.4. Informacijos pateikimas

Šioje dalyje galima peržvelgti analizuojamo teksto: kalbos dalis, išdetalizuotas žodžių kalbos dalis, atskirtus sakinius bei veiksmažodžių galūnes, tiesiog nuspaudus ant jų pasirinkimo mygtuko.



Tarkime pasirinkus „Kalbos dalys“ išvystame langą:

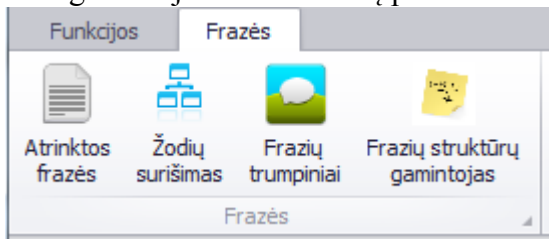
| Zodis | Kalbos Dalis | Zod Pasik Sk | Lemma | Lemma Pasik Sk | Eil Nr |
|----------|----------------------------|--------------|-----------------------|----------------|--------|
| Keturi | sktv., kiek., vyr. g., V | | keturi | | 1 |
| metai | dkt., vyr. g., dgs., V | | metai | | 11 |
| darbo | dkt., vyr. g., vns., K | | darbas | | 54 |
| triskart | prv., teig., nelygin. I | | triskart | | 13 |
| keitėsis | dlv., teig., sngr., vei... | | keistis(-čiasi,-tėsi) | | 15 |

44 Pav. Kalbos dalių langas

- Žodis – tekste esantis žodis.
- Kalbos dalis – rasto žodžio kalbos dalis
- Lema – žodžio veiksmažodinė forma
- Lemos pasikartojimo skaičius – Tekste rastų vienodų lemos pasikartojimo skaičius.
- Eil Nr – Žodžio eilės numeris

5.1.3. Frazės

Iš meniu pasirinkę „Frazės”, gauname sąrašą įrankių skirtų dirbti su frazėmis. Pirmi trys įrankiai (Atrinktos frazės, Žodžių surišimas, Frazės-sinonimai) skirti darbui su surastomis frazėmis. Frazijų struktūros gamintojas skirtas frazių paieškai.

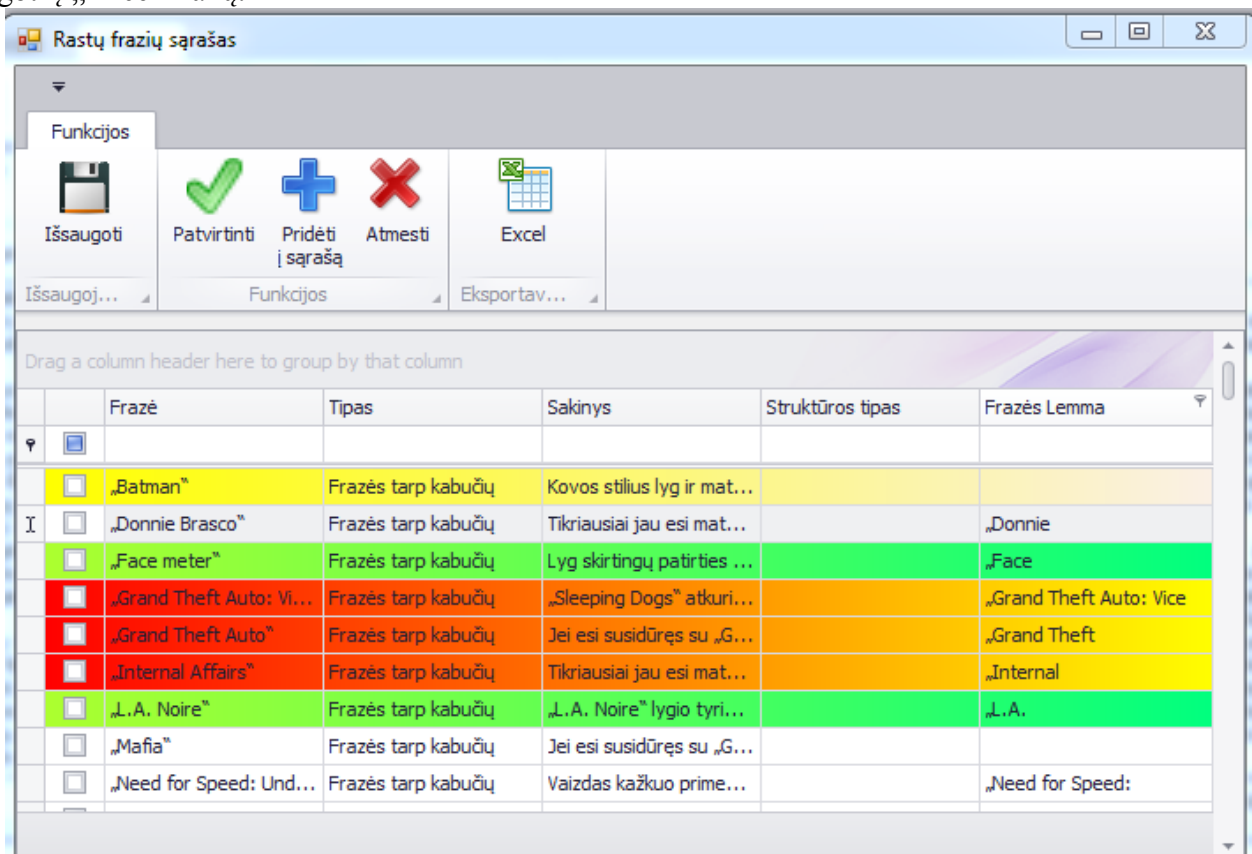


45 Pav. Frazijų meniu lango pasirinkimai

- Atrinktos frazės – atvaizduoti rastas frazes.
- Žodžių surišimas – leidžia padaryti faktus pvz. „Sleeping dog yra žaidimas“
- Frazė sinonimai – atvaizduoja rastu sinonimus pvz. KTU – Kauno Technologijos Universitetas
- Frazijų struktūrų gamintojas – leidžia sudaryti frazių šablonus, juos redaguoti .

5.1.3.1. Atrinktos frazės

Atrinktose frazėse matome informaciją apie fražę (frazės tipas, sakinytis, struktūros tipas, kuri rado fražę, frazės lemos). Yra penkios galimybės: išsaugoti, patvirtinti, pridėti į sąrašą, at mesti arba išsaugoti į „Excel“ failą.



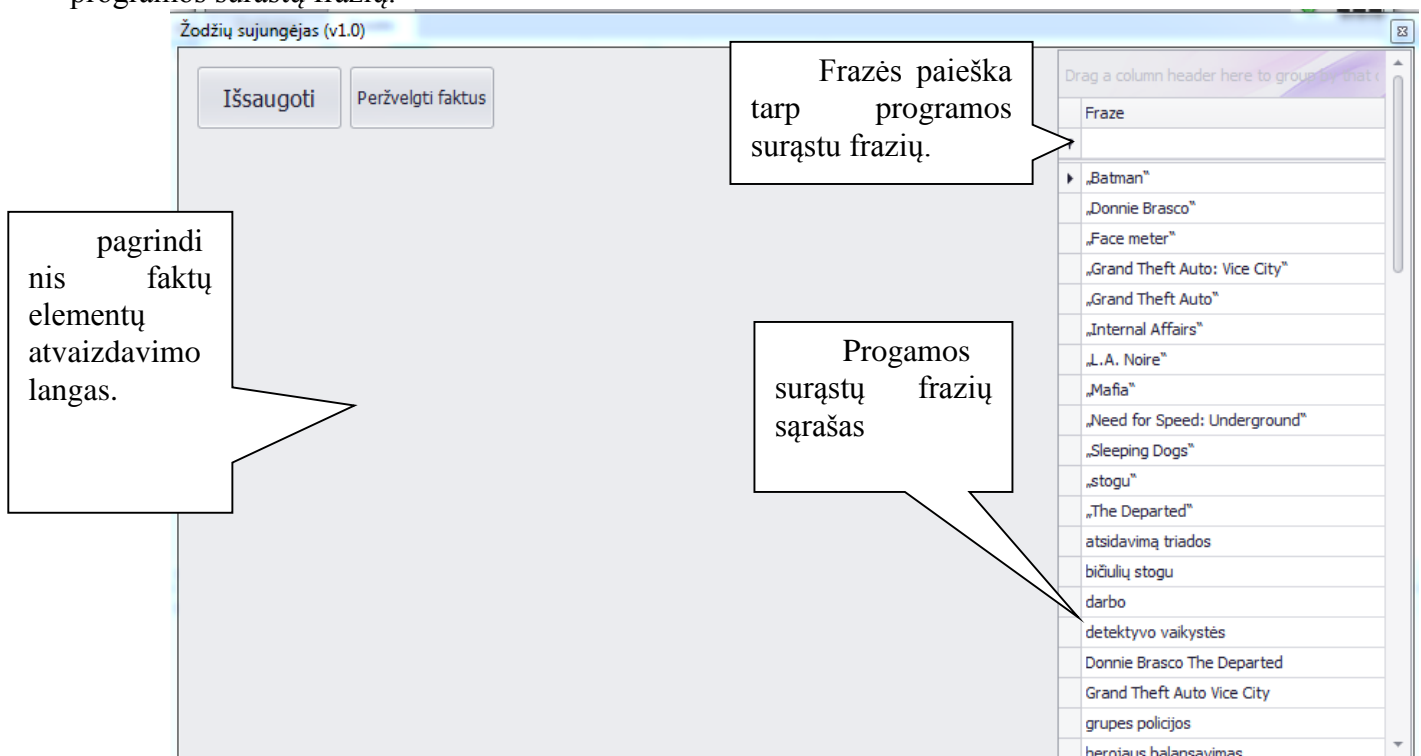
46 Pav. Atrinktų frazių pagrindinis langas

- Išsaugoti – išsaugo patvirtintas frazes ir panaikina atmetas frazes iš rastų frazių sąrašo.
- Patvirtinti – patvirtinti, kad programos surasta fražė yra tinkama. Patvirtintos frazės atvaizduojamos frazių sąrašė žalia spalva.

- Pridėti į sąrašą – galima į programos surastų frazių sąrašą pridėti naują frazę. Frazių sąrašė pridėtos frazės atvaizduojamos geltona spalva.
- Atmesti – leidžia atmesti surastas frazes kaip netinkamas. Frazių sąrašė jos atvaizduojamos raudona spalva.

5.1.3.2. Žodžių surišimas

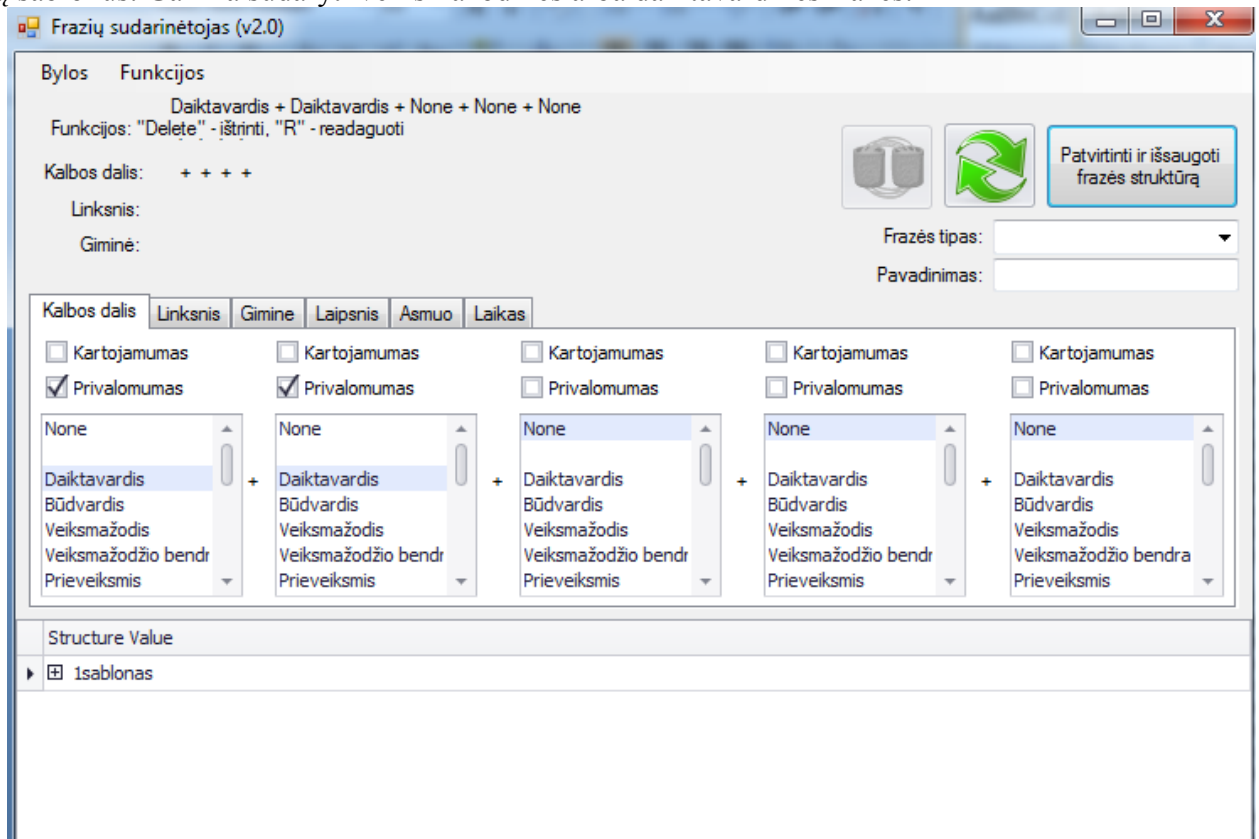
Galima rastas frazes tarpusavyje sujungti ir padaryti kaip faktą pvz. „Sleeping dog yra žaidimas.“. Faktą galima sudaryti iš programos surastų frazių arba įvesti savo. Faktai sudaryti iš elementų. Faktas elementas yra įrašas, tiek iš frazių sąrašo, tiek įvestas. Yra frazės paieška tarp programos surastų frazių.



47 Pav. Žodžių sujungėjo langas

5.1.3.3. Frazių šablonų gamintojas

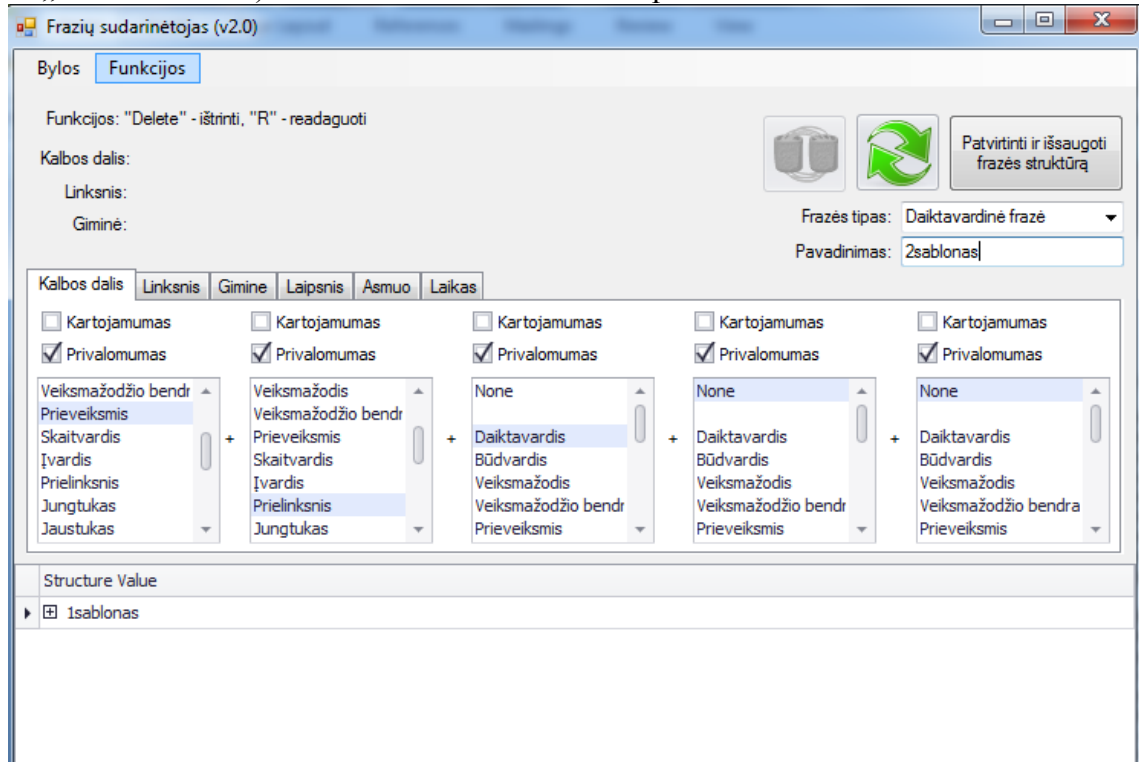
Žemiau yra pavaizduotas frazių šablono sudarytojo langas. Šiame lange yra galimybė sudaryti frazių šablonus. Galima sudaryti veiksmažodines arba daiktavardines frazes.



48 Pav. Frazių sudarinėtojo langas

5.1.3.3.1. Naujo šablono įvedimas

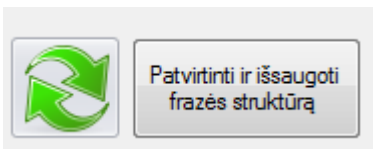
Galima pasirinkti, kad šablone būtų nurodyta kalbos dalis, linksnis, giminė, laipsnis, asmuo, laikas. Programa automatiškai atmets neleistinus pasirinkimus pvz. pasirinkus pirmame elemente būdvardis. Šiam elementui nebus leista pasirinkti laipsnio, asmens ir laiko. Iš viso gali būti nurodyti penki elementai. Tai pat gali būti nurodyta ar elementas daug kartų pasikartoja (šablonas yra būdvardis ir pažymėta „Kartojamumas“, sekantis elementas yra daiktavardis, sakinytis yra „Rudas , piktas šuo bėga“, tai šablonas ras „Rudas piktas šuo“). Yra galimybė pažymėti privalomumą (Šablonas sudarytas iš trijų elementų privalomas, prielinksnis neprivalomas ir daiktavardis privalomas. Du sakiniai „Gražiai į namus parėjau“ ir „Gražiai namas atrodo“. Šis šablonas ras dvi frazes „Gražiai į namus“ ir „Gražiai namas“). Turi būti būtinai nors vienas privalomas elementas.



49 Pav. Frazių sudarinėtojo langas

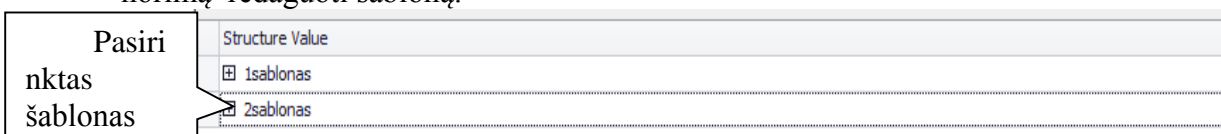
Norėdami išsaugoti šabloną reikia paspausti mygtuką „Patvirtinti ir išsaugoti frazės struktūrą“.

Skirta dar kartą išnagrinėti įkeltą tekstą su esančiomis ir naujai įvestais šablonais.



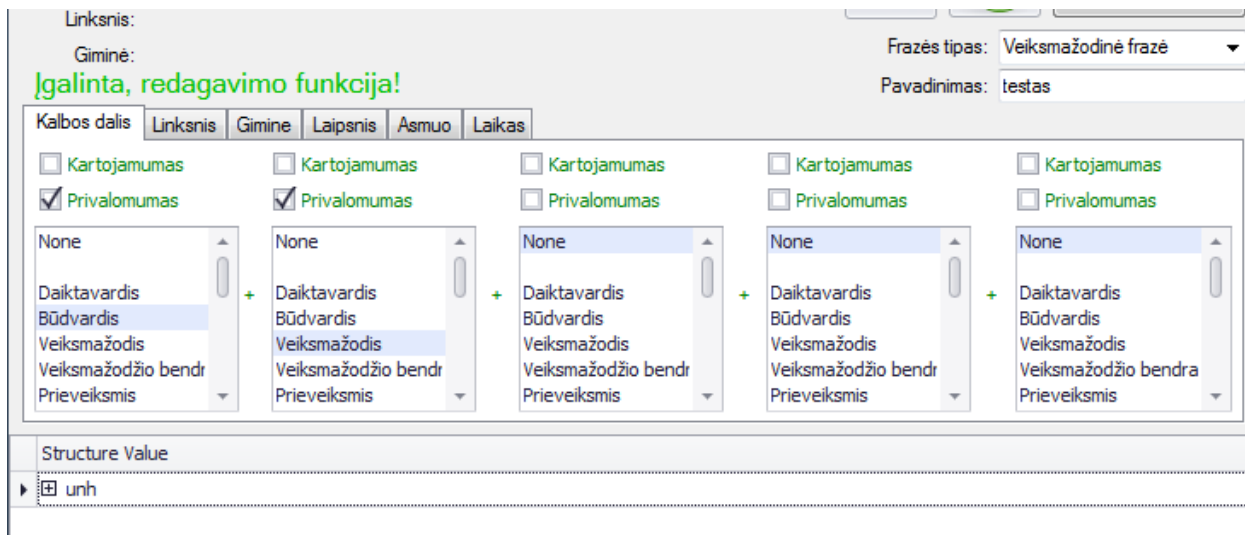
5.1.3.3.2. Šablono redagavimas

Yra galimybė išsaugotus šablonus redaguoti. Reikia iš išsaugotų šablonų sąrašo pasirinkti norimą redaguoti šabloną.



50 Pav. Frazių šablonų sąrašas

Pasirinkus norimą šabloną, reikia paspausti klaviatūros klavišą „R“

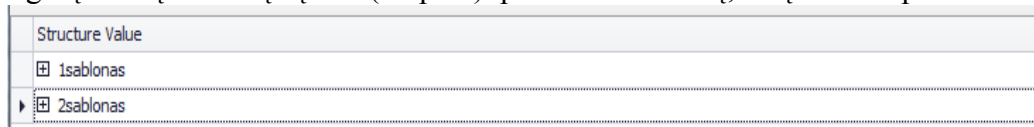


51 Pav. Frazių sudarinėtojo langas redagavimo režime

Jeigu esate redagavimo režime atsiranda užrašas „Įgalinta, redagavimo funkcija!“. Tada galima keisti šablono elementų reikšmes. Baigus keisti, kad išsaugotų pakeitimus reikia paspausti mygtuką „Patvirtinti ir išsaugoti frazės struktūrą“.

5.1.3.3.3. Šablono ištrynimasis

Yra galimybė nebereikalingą šabloną pašalinti, kad pagal jį būtų nebeieškoma frazių. Reikia iš išsaugotų frazių šablonų sąrašo (87 pav.) pasirinkti šabloną, kurį norima pašalinti



52 Pav. Frazių šablonų sąrašas

Reikia paspausti klaviatūros mygtuką „Delete“. Pažymėtas frazės šablonas yra pašalinamas.

5.2. Testavimo modelis

Pagrindinis testavimo tikslas yra pateikti ir įrodyti, jog parašyta programa veikia tinkamai ir kode ar klasėse ar metoduose nėra klaidų. Mažai šansų, jog ištestavus programą galima sakyti, jog programa nuo šiol veiks tinkamai ir gerai, tačiau testavimas sumažina blogo kodo riziką bei padidina programos efektyvumą bei lankstumą. Taipogi sistemos eksploatavimo metu kūrėjai įsipareigoja taisyti programos klaidas.

Remiantis testavimo planu pagrindiniai testavimo objektai bus klasės bei juose esantys metodai. Taipogi nemaža dalis testavimo apims ryšių testavimui tarp klasių bei metodų. Taipogi atliekant teminės srities žodyno testavimą, bus siekiama patikrinti ir užtikrinti, kad sistema atitinka sistemos specifikaciją.

5.2.1. Testavimo apimtis

Teminės srities žodyno sudarinėtojo programą galima suskaidyti į tokius elementus:

- Vartotojo sąsaja
- Sistemos moduliai
- Stambesnieji sistemos moduliai (tai sistemos moduliai kurie yra sudaryti iš kelėtos ar daugiau smulkesnių sistemos modulių)

Sistemos testavimui bus taikomi šie testavimo metodai: vienetų testavimas, integracijos testavimas ir aukšto lygio testavimas. Toliau pateiksime kurie sistemos elementai bus testuojami kiekvieno testavimo etapo metu:

Vienetų testavimas:

- Sistemos moduliai
 - Integruoti sistemos moduliai
- Šiame etape naudojamas juodos dėžės testavimas.

Integracijos testavimas:

- Vartotojo sąsaja
- Sistemos moduliai
- Integruoti sistemos moduliai

Aukšto lygio testavimo:

- Vartotojo sąsaja
- Sistemos moduliai

5.2.2. Testavimo strategija

5.2.2.1. Vienetų testavimas

Vienetų testavimo etape bus testuojami atskiri teminės srities žodyno sudarinėtojo moduliai ir paketai. Bus taikomas baltos dėžės testavimo metodas, kai testuotojas gali disponuoti programos kodu. Pirmiausiai bus išsiaiškinta, kokius duomenis reikia paduoti į atitinkamą programos modulį ir bus nustatomi laukiami rezultatai. Tada automatizuotai bus paduodami testiniai duomenys ir tikrinami rezultatai. Testavimo rezultatų matysime testavimo įrankio išvedimo lange, kuriame bus parašytas atitinkamas rezultatas, ar šis modulis praėjo testą ar ne.

Vienetų testavimas tai pats pirmas testavimo etapas, kurį atlieką patys programuotojai. Jie turi išsiaiškinti ar modulis veikia tinkamai. Jai modulis veikia, ne tinkamai, tai programuotojai privalo ištaisyti šias klaidas kuo skubiai, nes šio testavimo metu atsiradusios klaidos, gali pakenkti tolimesniam sistemos veikimui.

5.2.2.2. Aukštesnio lygmens testavimas

Šiame etape bus testuojami nefunkciniai sistemos reikalavimai. Teminės srities žodyno sudarinėtojo vienas svarbiausių nefunkcinių reikalavimų yra efektyvus kompiuterio resursų panaudojimas, taipogi lengva bei taipogi intuityvi bei turinti pakankamai funkcionalumo vartotojo sąsaja. Greitaveikos reikalavimai yra aprašyti sistemos specifikacijoje.

5.2.3. Testavimo resursai

Programiniai testavimo resursai:

- Teminės srities žodyno sudarinėtojas
- *UnitTest* įrankis, testams realizuoti

Žmogiškieji resursai:

- Programuotojas/testuotojas

5.3. Testavimo duomenys ir rezultatai

Pagrindiniai sistemos testavimo elementai:

- Vardų ir pavardžių atpažinimo algoritmas;
- Sutrumpinimų radimo algoritmas;
- Sakinių išskaidymo algoritmas;
- Frazijų radimo pagal šablonus (daiktvardinių frazių, veiksmažodinių frazių) algoritmai;

5.3.1. Vardų ir pavardžių algoritmo testavimas

Šio algoritmo testavimui buvo pasirinkti penki sakiniai kuriuose yra įvairaus tipo vardų bei pavardžių pvz.: Eimantas Žlabys, E.Žlabys, E.Ž., Valdas Mykolas Brazauskas ir t.t.

Žemiau pateiktuose lentelėse matyti gauti rezultatai iš įvairaus tipo sakinių.

- 1) *Kaip direktorių yra sake V.Adamkus man reikia kad perimtu E.Žlabys visa valdymą*
Gauti rezultatai: *V.Adamkus E.Žlabys*
- 2) *Kaip ir V.Adamkus, direktorių yra sake V.A. man reikia kad e.z. perimtu B.A. visa valdymą*
Gauti rezultatai: *V.Adamkus V.A. B.A.*
- 3) *Kaip direktorių yra sake V.A. man reikia kad e.z. perimtu B.A. visa valdymą*
Gauti rezultatai: *V.A. B.A.*
- 4) *Kaip direktorių yra sake Valdas Adamkus man reikia kad perimtu visa valdymą*
Gauti rezultatai: *Valdas Adamkus*
- 5) *Kaip direktorių yra sake Valdas Mykolas Brazauskas man reikia kad perimtu visa valdymą A.B.*
Gauti rezultatai: *Valdas Mykolas Brazauskas A.B.*

5.3.2. Sutrumpinimų radimo algoritmo testavimas

Šio algoritmo testavimui buvo pasirinkti įvairaus tipo sutrumpinimai kurie gali pretenduoti į sutrumpinimus, tačiau pridėti ir klaidingi elementai:

Sakinys: *KTU pats geriausias universitetas lyginant su VDU VGTU V ar Vu kitais Aukštaisiais LIETUVOS Universitetais*

Gauti rezultatai: *KTU VDU VGTU LIETUVOS*

5.3.3. Sakinių išskaidymo algoritmas

Šio algoritmo tikslas yra tinkamai išskirti sakinius iš teksto. Šio algoritmo testavimui buvo pasirinkti penki sakinių tipai kurie gali dažniausiai būti vartojami tekste. Taipogi prie kiekvieno sakinio prikabinamos įvairaus tipo galūnės kuriomis gali užsibaigti sakiniai (29 Lentelė Galimi sakinių pabaigos simboliai):

29 Lentelė Galimi sakinių pabaigos simboliai

| SAKINIO PABAIGOS ŽENKLAS | PAAIŠKNIMAS |
|--------------------------|-------------|
| . | Taškas |
| ? | Klaustukas |
| ! | Šauktukas |
| ... | Daugtaškis |

Sakiniai su kuriais buvo atliekami testavimai:

| |
|--|
| Šios sutarties sąlygos taikomos abonementui, fotogalerijai, (kartu visos paminėtos paslaugos vadinamos susijusiai programinei įrangai, žiniatinklio svetainei arba paslaugai (apskritai – paslaugoms) |
| Šios sutarties sąlygos taikomos „Microsoft! Hotmail“, „Microsoft.?!?!?! SkyDrive“, „Microsoft“ abonementui, „Windows Live Messenger“, „Windows“ fotogalerijai, „Windows Movie Maker“, „Microsoft Mail Desktop“, „Windows Live Writer“ (kartu visos paminėtos paslaugos vadinamos firminėmis „Microsoft“ paslaugomis), taip pat „Bing“, MSN, „Office.com“ ir bet kokiai kitai su šia sutartimi susijusiai programinei įrangai, žiniatinklio svetainei arba paslaugai (apskritai – paslaugoms) |
| Kaip direktoriu yra sake „ Man reikia kad V. Adamkus perimtu visa valdyma“ |
| 1. V. Adamkus visada daug padedavo saliai |
| 26. V. Adamkus teigė: „visada Ž. Eimantas bus prezidentas“ nuo 1989 metų |

5.3.4. Frazijų radimo pagal šablonus algoritmų testavimai

Šio algoritmo testavimui buvo pasirinktas vienas pagrindinis sakiny, kuriame yra įvairaus tipo kalbos dalių. Visi šablonai buvo generuojami automatiškai iš trijų elementų: Daiktavardžio, būdvardžio bei skaitvardžio. Iš šių trijų elementų buvo generuojamos nuo dviejų iki penkių ilgio šablonų elementai, kurie sudarydavo šabloną.

Testavimo elementai sudarė tokius atvejus:

- Visi šablono elementai yra daiktavardžiai;
- Šablono elementai maišosi tarpusavyje;

Veiksmožodinėms frazėms pagal šablonus testuoti buvo pasirinktas sakiny:

Studentai labai skubėjo į gražų Kauno Technologijos Universitetą.

Šį sakinį sudaro dvi dalys t.y. daiktvardinė bei veiksmožodinė frazė. Testavimui buvo pasirinkti keli tipai daiktvardinių bei veiksmožodinių frazių šablonai:

Veiksmožodinei frazei rasti: [Prieveiksmis] + Veiksmožodis + [Prielinksnis] + DF

Daiktavardinei frazei rasti: [Būdvardis] + Daiktavardis*

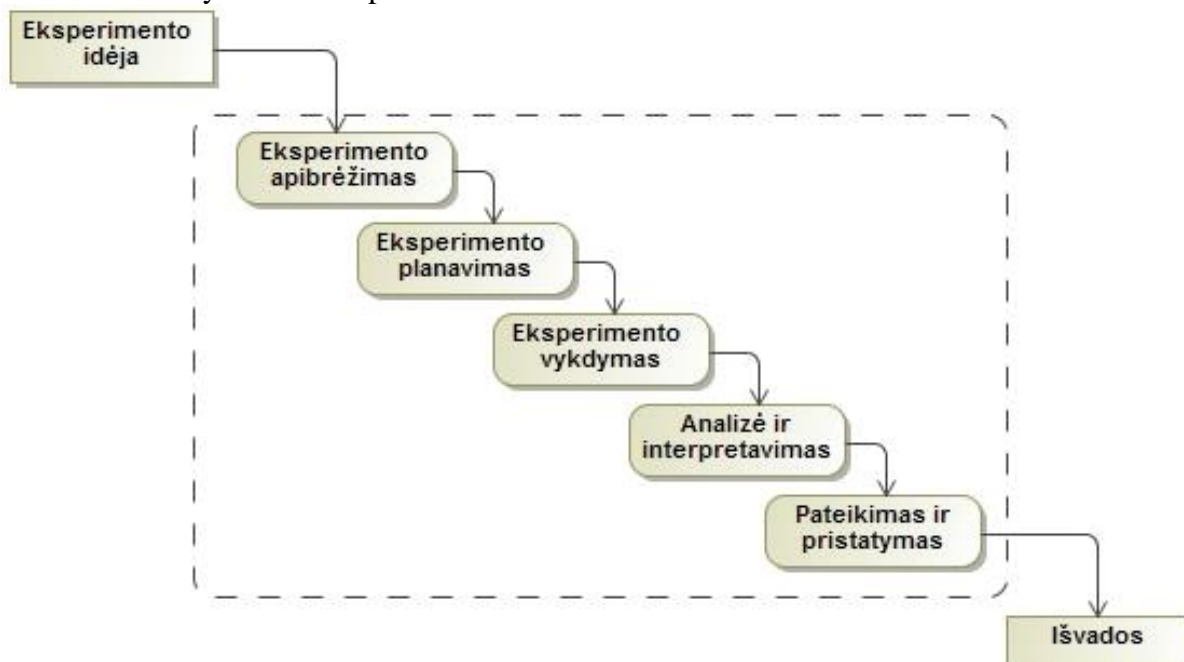
Gauti rezultatai:

Sistema išanalizavus sakinį pateikia šiuos rezultatus:

| Rezultatas | Tipas |
|---|-------|
| <i>gražų Kauno Technologijos Universitetą</i> | DF |
| <i>Kauno Technologijos Universitetą.</i> | DF |
| <i>skubėjo į</i> | VF |

6. EKSPERIMENTINIS AUTOMATINIO FRAZIŲ ATPAŽINIMO TEKSTE SUDARYTŲ ŠABLONŲ PAGRINDU SISTEMOS TYRIMAS

Žemiau (53 Pav. Eksperementavimo procesas) pateikiamas eksperimentavimo proceso schema kuria remiantis bus vykdomas eksperimentas.



53 Pav. Eksperementavimo procesas

6.1. Eksperimento planas

Eksperimento tikslas – pateikiant sistemai bei asmenims kurie dalyvauja bandyme tam tikrus tekstus, patikrinti programos algoritmo veikimą.

Eksperimentą atlieka eilinis vartotojas kuris išanalizuoja tekstus rankiniu būdu, po to tekstus duoda apdoroti programai. Taipogi eksperimente dalyvauja ekspertas kuris taipogi tuos tekstus išanalizuoja rankiniu būdu.

Eksperimento grėsmės:

- Dėl taisyklių, kurių pagalba algoritmas aptinka frazes iš vartotojo sudarytų šablonų neatitikimo.
- Esant netiksliam šablonui galimi nekorektiški rezultatai.
- Dėl temos naujumo informacinių sistemų srityje ir informacijos stokos algoritmo eksperimento rezultatai nebus visiškai korektiški.

Eksperimente naudojami įrankiai:

- Algoritmo sugebandžio naudojant šablonus tekste aptikti daiktavardines bei veiksmožodines frazes prototipas

Eksperimente naudojami duomenys:

- Tam tikri atrinkti tekstai.
- Atrinkti bei sudaryti šablonų tipai, kurie bus naudojami sistemai aptikti frazes

6.2. Eksperimento rezultatai

Šiame skyriuje pateikiami eksperimento rezultatai kurie gaunami naudojant frazių pagal šablonus atpažinimo algoritmą ir taipogi asmens, kuris atlieka eksperimentą rankiniu būdu rezultatai.

Eksperimentui atlikti sistemoje naudojamų šablonų sudėtis ir elementų kiekiai yra nekintantys. Eksperimentui atlikti sistemoje pasirinkti penki šablonai, kurie galimi dažniausiai atspindi frazę.

Žemiau pateiktoje (30 Lentelė Sistemos ekeperimente naudojamų šablonų sąrašas) lentelėje išrašyti visų šablonų sudedamosios dalys:

30 Lentelė Sistemos ekeperimente naudojamų šablonų sąrašas

| <i>Šablonas</i> | <i>Tipas</i> | <i>Paaiškinimas</i> |
|----------------------|---------------------|---|
| $[Sk][B^*][D^*](k)D$ | Daiktvardinė frazė | Šabloną sudaro skaitvardis neprivalomas, būdvardis, kuris gali kartotis tačiau jis neprivalomas, daiktavardis kilmininko linksnyje, kuris gali kartotis tačiau neprivalomas ir daiktavardis |
| $Dal[B^*][D^*](k)D$ | Daiktvardinė frazė | Šabloną sudaro dalyvis, būdvardis, kuris gali kartotis tačiau jis neprivalomas, daiktavardis kilmininko linksnyje, kuris gali kartotis tačiau neprivalomas ir daiktavardis |
| $DVbendr$ | Daiktvardinė frazė | Šabloną sudaro daiktavardis ir veiksmažodžio bendratis |
| $V[Prl]$ | Daiktvardinė frazė | Šabloną sudaro veiksmažodis ir prielinksnis, kuris yra neprivalomas |
| $VVbenr[Prl]$ | Veiksmažodinė frazė | Šabloną sudaro veiksmažodis, veiksmažodžio bendratis ir prielinksnis, kuris yra neprivalomas |

Kadangi sistema automatiškai, nepriklausomai ar šablonai yra sudaryti ar ne ji ieško tekste vardų ir pavardžių, sutrumpinimų. Taigi prie šio eksperimento prisidės ne tik šablonų rastos frazės tačiau ir automatiškai tekste rastos frazės.

Algoritmui bei asmeniui atliekančiam eksperimentą yra pateikiami penki tekstai kurie yra bendros tematikos – politika. Šie tekstai yra patalpinti internete adresu: <https://eimzlab.stud.if.ktu.lt/Tekstai%20analizei/>.

Teksto analizei bei frazių radimui buvo išskelti kriterijai:

- Veikimo tikslumas (palyginimas sistemos ir asmens atliekančio eksperimentą)
- Kiek frazių buvo rasta tekste
- Kiek frazių yra korektiškų kiek nekorektiškų
- Asmens atlikusio eksperimentą bei sistemos rastų frazių palyginimas

6.2.1. Sukurtos sistemos eksperimento rezultatai

Sistemos algoritmas tekstuose surado 296 frazes. Žemiau esančioje (31 Lentelė Sistemos kiekvieno šablono rastų frazių skaičius) lentelėje yra pateikiamos tikslesnės kiekvieno šablono aptiktų frazių skaičius:

31 Lentelė Sistemos kiekvieno šablono rastų frazių skaičius

| <i>Šablono pavadinimas</i> | <i>Rastų frazių skaičius</i> |
|----------------------------|------------------------------|
| $[Sk][B^*][D^*](k)D$ | 151 |
| $Dal[B^*][D^*](k)D$ | 43 |
| $DVbendr$ | 10 |
| $V[Prl]$ | 62 |
| $VVbenr[Prl]$ | 6 |
| <i>Viso:</i> | 272 |

Kaip matyti iš lentelės frazių atpažinimo naudojant šablonus algoritmas aptiko 238 frazes kurios tenkino sudarytus šablonus. Kitas 24 frazes sistema aptiko automatiškai ir juos priskyrė vardams bei pavardėms ir sutrumpinimams:

32 Lentelė Ne šablonų rastų razių skaičius

| <i>Šablono pavadinimas</i> | <i>Rastų frazių skaičius</i> |
|----------------------------|------------------------------|
| <i>Sutrumpintos frazės</i> | 1 |
| <i>Vardai bei pavardės</i> | 23 |
| <i>Viso:</i> | 24 |

Žemiau esančioje (33 Lentelė Sistemos rastų frazių kiekiai kiekviename tekste) lentelėje yra pateikiamas frazių skaičius rastas tam tikrame tekste. Taipogi 1priede pateikiamos visos sistemos rastos frazės:

33 Lentelė Sistemos rastų frazių kiekiai kiekviename tekste

| <i>Teksto pavadinimas</i> | <i>Rastų frazių skaičius</i> |
|------------------------------------|------------------------------|
| <i>bals-Pasaulis-090907-2.txt</i> | 61 |
| <i>bals-Lietuva-091220-111.txt</i> | 27 |
| <i>bals-Lietuva-091213-104.txt</i> | 48 |
| <i>bals-Lietuva-091203-94.txt</i> | 81 |
| <i>bals-Lietuva-091124-85.txt</i> | 79 |
| <i>Viso:</i> | 296 |

6.2.2. Eilinio vartotojo atlikusio eksperimentą rezultatai

Asmuo atlikęs teksto analizę yra eilinis vartotojas, kuris sugeba tekste aptikti frazes, taipgi žino pagrindines kalbos dalis, tekstuose surado 74 frazes. Žemiau esančioje (34 Lentelė Eilinio vartotojo rastų frazių skaičius kiekviename tekste) lentelėje yra pateikiama kiek eilinis vartotojas aptiko frazių kiekviename tekste, vartotojo rastos frazės pateikiamos (2 priede):

34 Lentelė Eilinio vartotojo rastų frazių skaičius kiekviename tekste

| <i>Teksto pavadinimas</i> | <i>Rastų frazių skaičius</i> |
|------------------------------------|------------------------------|
| <i>bals-Pasaulis-090907-2.txt</i> | 20 |
| <i>bals-Lietuva-091220-111.txt</i> | 11 |
| <i>bals-Lietuva-091213-104.txt</i> | 19 |
| <i>bals-Lietuva-091203-94.txt</i> | 14 |
| <i>bals-Lietuva-091124-85.txt</i> | 10 |
| <i>Viso:</i> | 74 |

6.2.3. Eksperto atlikusio eksperimentą rezultatai

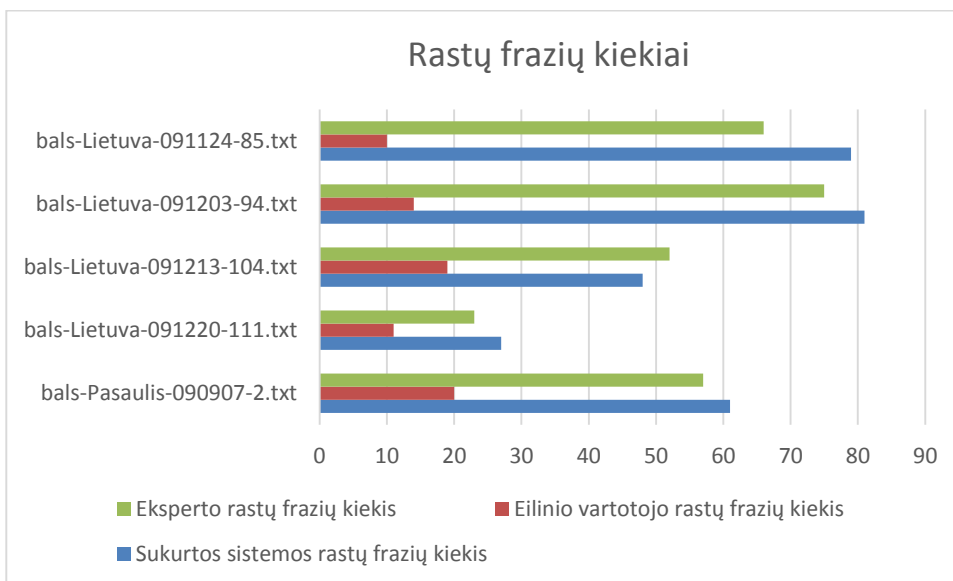
Asmuo atlikęs teksto analizę yra ekspertas, kuris gerai žino Lietuvių kalbos tekstų struktūra. Suvokia bei sugeba sudaryti bei rasti frazes. Taipogi asmuo sugeba ir taiko gramatines kalbos dalis frazėms aptikti. Ekspertas tekstuose surado 293 frazių. Žemiau esančioje (35 Lentelė Eksperto rastų frazių skaičius kiekviename tekste) lentelėje yra pateikiamos kiek kiekviename tekste ekspertas patiko frazių:

35 Lentelė Eksperto rastų frazių skaičius kiekviename tekste

| <i>Teksto pavadinimas</i> | <i>Rastų frazių skaičius</i> |
|------------------------------------|------------------------------|
| <i>bals-Pasaulis-090907-2.txt</i> | 77 |
| <i>bals-Lietuva-091220-111.txt</i> | 23 |
| <i>bals-Lietuva-091213-104.txt</i> | 52 |
| <i>bals-Lietuva-091203-94.txt</i> | 75 |
| <i>bals-Lietuva-091124-85.txt</i> | 66 |
| <i>Viso:</i> | 293 |

6.3. Sprendimo savybių analizė, kokybės kriterijų įvertinimas

Atlikus eksperimentus iš sistemos pusės bei asmenų kurie atliko teksto analizę, galime matyti, jog rastų frazių skaičius skiriasi. Žemiau pateiktoje diagramoje (54 Pav. Rastų frazių kiekiai) matyti kaip skiriasi frazių skaičius tam tikrame tekste.



54 Pav. Rastų frazių kiekiai

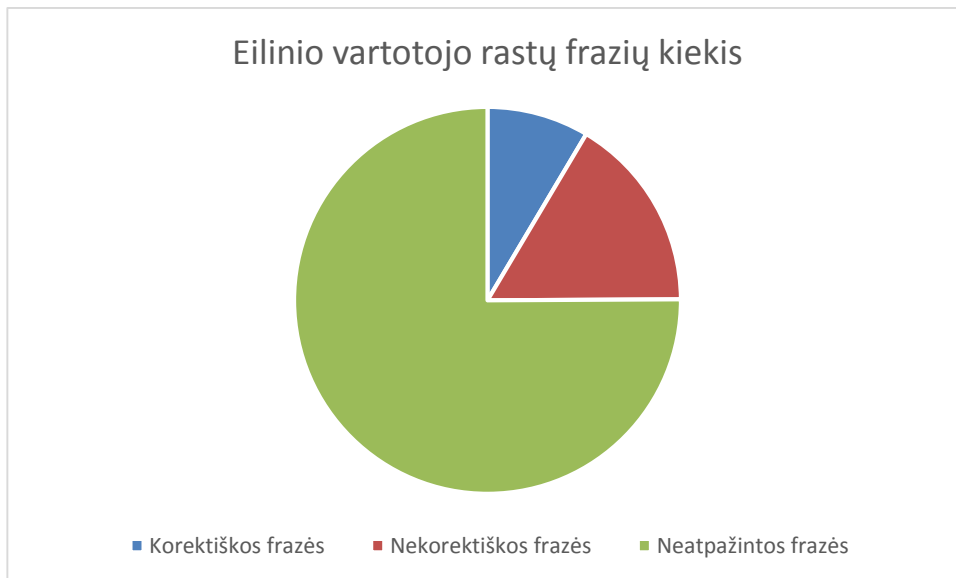
Iš viso rastų frazių kiekiai pateikti (55 Pav. Visuose tekstuose rastų frazių kiekiai) diagramoje:



55 Pav. Visuose tekstuose rastų frazių kiekiai

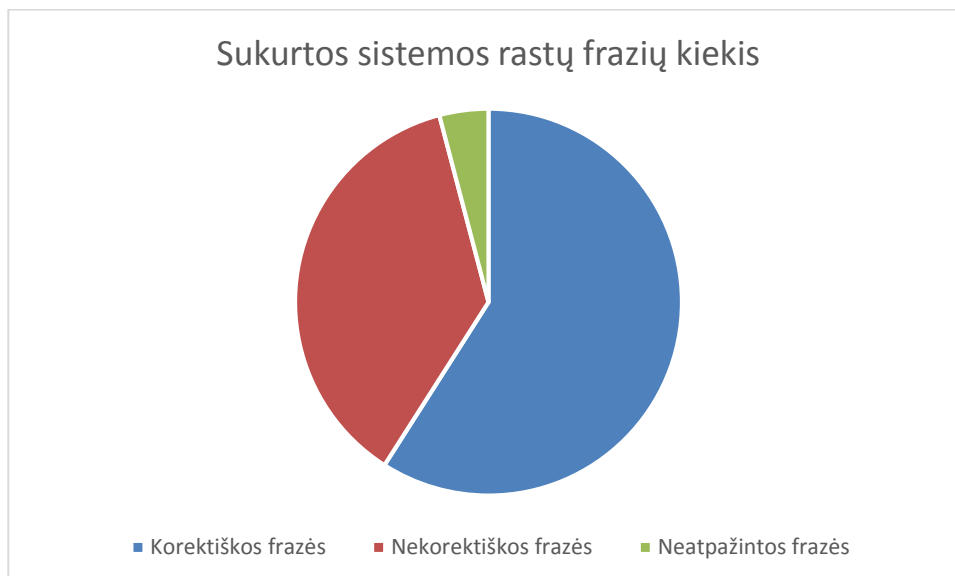
Žemiau pateiktose diagramose pavaizduoti kiek sukurtos sistemos (57 Pav. Sukurtos sistemos rastų frazių palyginimas su eksperto rastomis frazėmis) ir eilinio vartotojo (56 Pav. Eilinio vartotojo rastų frazių palyginimas su eksperto rastomis frazėmis) lyginant su ekspertu rastų frazių skirtumai. Lyginant rastas frazes buvo iškelti trys kriterijai, pagal kuriuos buvo matuojama:

- Korektiškos frazės – frazės rastos tiek eksperto tiek sukurtos sistemos yra vienodos;
- Nekorektiškos frazės – frazės rastos sukurtos sistemos nepilnai ar dalinai atitinka frazes rastos eksperto;
- Neatpažintos frazės – frazės rado ekspertas, tačiau sukurta sistema jų nerado;



56 Pav. Eilinio vartotojo rastų frazių palyginimas su eksperto rastomis frazėmis

Kaip matyti iš 54pav. jog vartotojo rastų frazių skaičius lyginant su eksperto ganėtinai skiriasi. Vartotojas korektiškai rado tik 25 frazes lyginant su eksperto rastomis frazėmis. Nekorektiškai rastos, ar iš dalies panašios frazės su eksperto sudaro 48 frazes. Tuo tarpu nerastų frazių skaičius siekia 220, t.y. vartotojas neišskyrė frazių, kurias rado ekspertas.

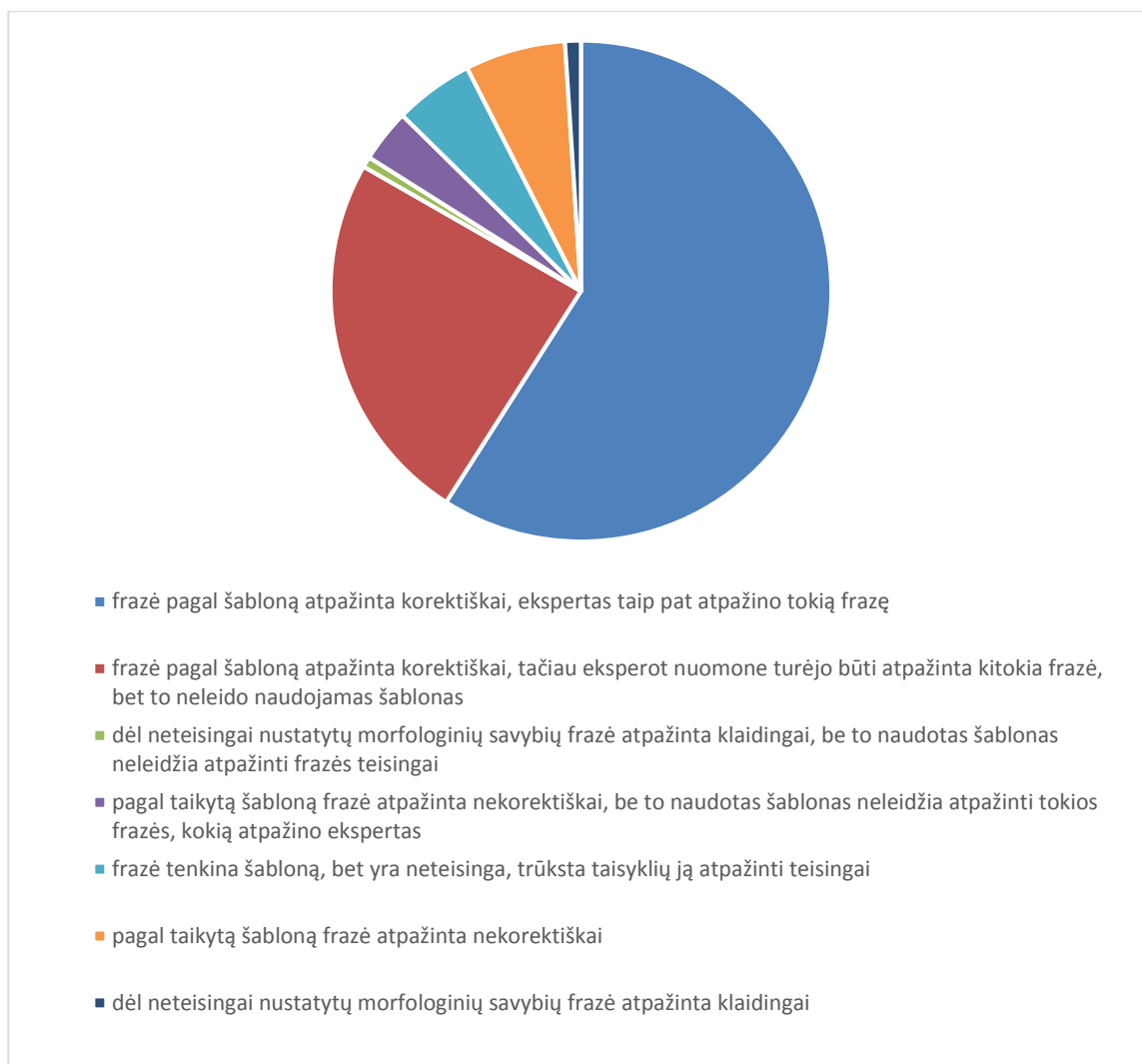


57 Pav. Sukurtos sistemos rastų frazių palyginimas su eksperto rastomis frazėmis

Iš 55pav. matyti jog sukurta sistema surado daugiau nei pusę, t.y. 173 frazes kurias surado ir ekspertas. Neatpažintų frazių pasitaikė tik 12, t.y. ekspertas rado tačiau sistema nerado. Nekorektiškai rastų frazių skaičius siekia 108, taip gali nutikti dėl nepilnų ar netinkamų šablonų kurie buvo sudaryti sistemai.

Žemiau pateiktoje diagramoje mes sukurtos sistemos rastas frazes išskaidėme į kitus septynis kriterijus, pagal kuriuos galime matyti sistemos rastų frazių tikslumą žiūrint pagal eksperto rastas frazes. Išskelti kriterijai:

- frazė pagal šabloną atpažinta korektiškai, ekspertas taip pat atpažino tokią frazę
- frazė pagal šabloną atpažinta korektiškai, tačiau eksperto nuomone turėjo būti atpažinta kitokia frazė, bet to neleido naudojamas šablonas
- dėl neteisingai nustatytų morfologinių savybių frazė atpažinta klaidingai, be to naudotas šablonas neleidžia atpažinti frazės teisingai
- pagal taikytą šabloną frazė atpažinta nekorektiškai, be to naudotas šablonas neleidžia atpažinti tokios frazės, kokią atpažino ekspertas
- frazė tenkina šabloną, bet yra neteisinga, trūksta taisyklių ją atpažinti teisingai
- pagal taikytą šabloną frazė atpažinta nekorektiškai
- dėl neteisingai nustatytų morfologinių savybių frazė atpažinta klaidingai



58 Pav. Sukurtos sistemos atpažintų frazių detalesnis pjūvis

36 Lentelė Sistemai keltų kriterijų sąrašas bei rezultatai

| Kriterijus | |
|---|-----|
| <i>frazė pagal šabloną atpažinta korektiškai, ekspertas taip pat atpažino tokią frazę</i> | 173 |
| <i>frazė pagal šabloną atpažinta korektiškai, tačiau eksperto nuomone turėjo būti atpažinta kitokia frazė, bet to neleido naudojamas šablonas</i> | 71 |
| <i>dėl neteisingai nustatytų morfologinių savybių frazė atpažinta klaidingai, be to naudotas šablonas neleidžia atpažinti frazės teisingai</i> | 2 |
| <i>pagal taikytą šabloną frazė atpažinta nekorektiškai, be to naudotas šablonas neleidžia atpažinti tokios frazės, kokią atpažino ekspertas</i> | 10 |
| <i>frazė tenkina šabloną, bet yra neteisinga, trūksta taisyklių ją atpažinti teisingai</i> | 15 |
| <i>pagal taikytą šabloną frazė atpažinta nekorektiškai</i> | 19 |
| <i>dėl neteisingai nustatytų morfologinių savybių frazė atpažinta klaidingai</i> | 3 |

Kaip matyti ir iš aukščiau esančios (36 Lentelė Sistemai keltų kriterijų sąrašas bei rezultatai) lentelės bei 56pav. jog sistemos rastos frazės ir eksperto rasto frazės daugiau nei pusė sutampa lyginant su visomis rastomis frazėmis. Pirmąjį kriterijų, t.y. *frazė pagal šabloną atpažinta korektiškai, ekspertas taip pat atpažino tokią frazę* atitiko 173 frazės. Toliau kitos sistemos frazės išsibarstė į kitus kriterijus.

Sistemos frazių radimo korektiškumą galima padidinti sudarant tinkamesnius šablonus. Sudarant šablonus būtina remtis lietuvių kalbos žiniomis, norint patikslinti šablono dedamąsias. Dėl neteisingų, ar nepilnų sistemos šablonų 71 frazė buvo rasta nepilna, lyginant eksperto. Tokių rastų frazių pavyzdžiai pateikti žemiau esančioje (37 Lentelė Nepilnų frazių lyginant su eksperto sąrašas) lentelėje:

37 Lentelė Nepilnų frazių lyginant su eksperto sąrašas

| <i>Sistemos rasta frazė</i> | <i>Eksperto rasta frazė</i> |
|--|---|
| <i>švietimo</i> | Šiaulių švietimo darbuotojų |
| <i>žadama</i> | žadama išmokėti |
| <i>paramos ir labdaros fondo Kalėdų šventėje</i> | surengtoje Almos Adamkienės paramos ir labdaros fondo Kalėdų šventėje |
| <i>pirmoji šalies ponis</i> | buvusi pirmoji šalies ponis |
| <i>prezidento</i> | prezidento Valdo Adamkaus žmona Alma |

Norint pagerinti automatinį sistemos frazių radimą, reikia sudaryti tinkamus bei gerus šablonus. Esant netinkamiems, ar ne pilniems šablonams, aptikti reikiamas frazes sudėtinga. Norint sudaryti gerus šablonus reikia gerai išmanyti Lietuvių kalbos gramatiką, žinoti pagrindines daiktavardinių bei veiksmažodinių frazių sudarinėjimo taisykles.

Taipogi sukurtos sistemos nepilnų frazių radimo aspektas yra nepilnas sistemos išbaigtumas. Algoritmų nepilnumas, ne visų taisyklių įgyvendinimas. Aptikti sudėtingesnės struktūros frazės yra ypač sudėtinga, kadangi sunku sudaryti jiems šablonus.

7. IŠVADOS

- 1) Atlikus panašių sistemų analizę, buvo pastebėta, jog nei viena iš analizuotų sistemų netinka Lietuviško teksto analizei. Taipogi nei viena iš sistemų neturėjo galimybės pačiam vartotojui nustatyti ieškomų frazių, junginių ar kitokių struktūrų žodžių, sudaryti numanomas šablonus pagal kuriuos sistema galėtų rasti bei aptikti frazes ar junginius.
- 2) Atlikta analizė leidžia teigti, jog frazių ar žodžių junginių radimas tekstuose ar žodynų sudarinėjimas yra gan sudėtingas ir laiko reikalaujantis mechanizmas.
- 3) Iš atliktų teminės srities žodynų sudarinėtojų programų analizės galima teigti, jog Lietuviško teksto žodynų sudarymas gali būti gan sudėtingas. Norėdami palengvinti frazių radimą bei pačios programos kūrimą galima bus remtis, jau esamomis kitoms kalboms skirtomis funkcijomis bei paieškos metodais.
- 4) Nustačius pagrindinius tikslus, algoritmo funkcionalumus ir prototipo viziją buvo nustatytos pagrindinės algoritmo veiklos detalės, suprojektuotas ir atvaizduotas algoritmo funkcionalumas grafiškai bei aprašyti algoritmo veikimo principai.
- 5) Projektavimo metu buvo apibrėžti pagrindiniai algoritmo prototipo reikalavimai. Sukurtos taisyklės pagal kurias algoritmas vykdo numatytus funkcionalumus. Taip pat aprašyti bei atvaizduoti algoritmo veikimo žingsniai kuriais remiantis algoritmas pateikia rezultatus.
- 6) Atlikus eksperimentą buvo nustatyta jog sistema gali aptikti daugiau nei pusę frazių lyginant su eksperto rastomis. Eksperimento metu buvo prieita išvados jog norint geriau ir tinkamiau rasti frazes, reikia geriau sudaryti šablonus, atitinkamai parinkti tinkamesnius šablono elementus.
- 7) Remiantis darbo analizės, projektavimo ir algoritmo aprašymo medžiaga galime teigti, jog turime visapusišką įrankio prototipą, kuris gali pagreitinti bei padėti lengviau ir tiksliau rasti daiktvardines ir veiksmažodines frazes lietuviškame tekste. Algoritmas žymiai paspartina bei pagerina SBVR žodynų sudarymo eigą, taipogi pagerintų šių žodynų kokybę bei tikslumą.
- 8) Programos realizavimas parodė, jog frazių paieškos algoritmą naudojant šablonus įmanoma sukurti, o testavimas - , kad sukurtas bei realizuotas algoritmas veikia tinkamai, nes tenkina išskeltus reikalavimus.

8. LITERATŪRA

- [1] VDU, „<http://tekstynas.vdu.lt/>“, 1998-2013. [Tinkle]. Available: <http://tekstynas.vdu.lt/page.xhtml;jsessionid=5A6DA6C0EC9EAD7A1CAFD94713D3D9F8?id=morphological-annotator>. [Kreiptasi 13 balandžio 2014].
- [2] „dz.lki.lt“, 2012-2013. [Tinkle]. Available: <http://dz.lki.lt>. [Kreiptasi 09 balandžio 2014].
- [3] „Gerpikime žodį“, [Tinkle]. Available: http://ualgiman.dtiltas.lt/leksikografija_ir_zodynu_tipai.html. [Kreiptasi 09 balandžio 2014].
- [4] „Gerpikime žodį“, [Tinkle]. Available: http://ualgiman.dtiltas.lt/leksikografija_ir_zodynu_tipai.html. [Kreiptasi 09 balandžio 2014].
- [5] Labutis, Labutis, 1998.
- [6] SDL, „<http://www.translationzone.com/>“, SDL, 2011. [Tinkle]. Available: <http://www.translationzone.com/products/sdl-multiterm/desktop/>. [Kreiptasi 13 balandžio 2014].
- [7] Terminotix, „<http://www.terminotix.com/>“, Terminotix, [Tinkle]. Available: <http://www.terminotix.com/index.asp?name=SynchroTerm&content=item&brand=4&item=7&lang=en>. [Kreiptasi 13 balandžio 2014].
- [8] WordFast, „<http://www.wordfast.net/>“, Wordfast LLC, [Tinkle]. Available: <http://www.wordfast.net/index.php?whichpage=products&lang=engb>. [Kreiptasi 13 balandžio 2014].
- [9] R. J. B. a. B. K. B. Younga Park, „Automatic Glossary Extraction: Beyond Terminology Identification“, p. 7.
- [10] E. Vyšniauskas ir L. Nemuraitė, „Transforming Ontology Representation from OWL to Relational Database“, *Information Technology and Control*, t. 35A, nr. 3, p. 333–343, 2006.
- [11] K. Masiulis ir A. Krupavičius, Valstybės tarnyba Lietuvoje: praeitis ir dabartis: kolektyvinė monografija, Vilnius: Praction, 2007, p. 430.
- [12] V. Biržiška, „Spaudos draudimo klausimai“, *Kultūra*, nr. 5, pp. 249-235, 1929.
- [13] I. Valiulytė, „Išlaidos krašto apsaugai, jų pagrįstumas ir tikslingumas“, vasaris 2000. [Tinkle]. Available: <http://www.sociumas.lt>. [Kreiptasi 12 gruodžio 2001].
- [14] D. U. Library, „IEEE Citation style guide“, 2009. [Tinkle]. Available: http://libraries.dal.ca/content/dam/dalhousie/pdf/library/Style_Guides/IEEE_Citation_Style_Guide.pdf. [Kreiptasi 11 04 2013].
- [15] R. Gradauskas, „Hibridinis velomobilis“, įtraukta *Transporto priemonės - 99*, Kaunas, 2000.

9. PRIEDAI

9.1. 1 priedas. Sistemos rastos frazės

| Phrase | StructureName | TextName |
|-----------------------------------|------------------|-----------------------------|
| 10 ame | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 10 procentų | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| 12 | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| 13 | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 16 kvadratinių metrų kambarys | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 17 | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 250 | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| 41 metų | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| 58 metų vyras | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 6 kvadratinis metrus | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 6240 litų baudą | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| 76 metų namo savininkė | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| 7800 litų bauda | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| A.Adamkienei | Vardai pavardės | bals-Lietuva-091220-111.txt |
| A.Pukelio | Vardai pavardės | bals-Lietuva-091203-94.txt |
| A.Pukelis | Vardai pavardės | bals-Lietuva-091203-94.txt |
| aiškino | V[PrI] | bals-Lietuva-091124-85.txt |
| algų | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| Alma Adamkienė | Vardai pavardės | bals-Lietuva-091220-111.txt |
| Almos Adamkienės | Vardai pavardės | bals-Lietuva-091220-111.txt |
| AP. | Vardai pavardės | bals-Pasaulis-090907-2.txt |
| Apygardos teismą | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| Apygardos teismo teisėjų kolegija | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| apylinkės | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| apspjovęs | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| Arūnų Pukelį | Vardai pavardės | bals-Lietuva-091203-94.txt |
| asilu | DVBendr | bals-Lietuva-091203-94.txt |
| atlikta operacija | Dal[b*][D*](k)D | bals-Lietuva-091220-111.txt |
| atlyginimą | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| Atlyginimus | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| atostogų | DVBendr | bals-Lietuva-091213-104.txt |
| atsakomybės | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| atsisakė | V[PrI] | bals-Pasaulis-090907-2.txt |
| atsisakymą pateikti | DVBendr | bals-Pasaulis-090907-2.txt |
| atsiųstas ženklas | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| atstovų | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| atšventusio tauragiškio mašiną | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| atvira liepsna | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| Audi" Švinius | Vardai pavardės | bals-Lietuva-091203-94.txt |
| aukšto | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| avansus | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| baigęs prezidentas | Dal[b*][D*](k)D | bals-Lietuva-091220-111.txt |
| baltas | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| batus | DVBendr | bals-Pasaulis-090907-2.txt |

| | | |
|---------------------------------------|------------------|-----------------------------|
| <i>biudžetinėms įstaigoms</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>biurokratinių reikalavimų</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>būtina</i> | Dal[b*][D*](k)D | bals-Pasaulis-090907-2.txt |
| <i>buvę baldai</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>buvo</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>buvusi</i> | Dal[b*][D*](k)D | bals-Lietuva-091220-111.txt |
| <i>Buvusio kandidato teigimu</i> | Dal[b*][D*](k)D | bals-Pasaulis-090907-2.txt |
| <i>dabartiniu Niujorko meru</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>darbuotojos</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>darželio problema</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>degusio</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>delspinigių skaičiavimo</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>didesnių problemų</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>dirbs be</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>dolerių bauda</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>du metrus</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>dviejų nepilnamečių valdininko</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>vaikų akivaizdoje</i> | | |
| <i>finansinę padėtį</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>gaisrą</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Gaisrininkai</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Gaisro metu</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>gali pasiguosti</i> | VVbendr[prI] | bals-Lietuva-091203-94.txt |
| <i>gatvės</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Gatvės muzikantas</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>gavęs</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>gavo</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>gavo</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>gyvenimo</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>grasinimas sudeginti</i> | DVBendr | bals-Lietuva-091203-94.txt |
| <i>grasinimus</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>gruodį</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>guli</i> | V[PrI] | bals-Lietuva-091220-111.txt |
| <i>ir gėjumi</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>ir po</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>ir skrybėlę</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>ir visuomeninio transporto</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>naujovių įgyvendinimo</i> | | |
| <i>išdegė</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>išgelbėjo per</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>iškoneveikė</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>išnaudotas metų mokes fondas</i> | Dal[b*][D*](k)D | bals-Lietuva-091213-104.txt |
| <i>išpuolio viešoje</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>išteisintas</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>išvadino</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>išvakarėse</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>išvežtas</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>įspūdžio</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>įstaigų vadovėms</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |

| | | |
|------------------------------------|---------------------|-----------------------------|
| <i>įtakos vietos politikams</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>įteikė</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>yra</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>yra</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>yra apie</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>jaučiasi</i> | V[PrI] | bals-Lietuva-091220-111.txt |
| <i>JAV</i> | Sutrumpintos frazės | bals-Pasaulis-090907-2.txt |
| <i>jubiliejų</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>kadenciją</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>kaltinusį prokurorą</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>kambario</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>kambaryje</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>kandidatūrą</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>kartus</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>kaubojiškus batus</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>ketina konkuruoti su</i> | VVbendr[prI] | bals-Pasaulis-090907-2.txt |
| <i>ketinimų įgyvendinti</i> | DVBendr | bals-Lietuva-091124-85.txt |
| <i>kilęs gaisras</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>kilo nuo</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>konstatavo</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>kreipėsi</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>kreipėsi į</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>lapkritį</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>lašu anot</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>liepos mėnesį</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>liepsnas</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>ligoninėje</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>ligoninę</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>likimo</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>liko be</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>lyties</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>LNK Žinios.</i> | Vardai pavardės | bals-Lietuva-091220-111.txt |
| <i>Maiklu Blumbergu</i> | Vardai pavardės | bals-Pasaulis-090907-2.txt |
| <i>Manoma</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>mažesnėmis galimybėmis</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>mediniame gyvenamajame name</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Mercedes Benz" Tauragėje</i> | Vardai pavardės | bals-Lietuva-091203-94.txt |
| <i>mero postą</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>metu</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>metų birželio pradžioje</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>metų viduryje</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>miesto biudžete</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>Miesto darželių vedėjos</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>miesto vadovo rinkimuose</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>mirė</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>mokesčių</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>mokos fondus</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>namas</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |

| | | |
|------------------------------------|------------------|-----------------------------|
| <i>naujienu agentūra</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>neabejoja</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>nebeliko</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>nebesitiki</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>nebuvo</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>necenzūriniais žodžiais</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>nemokamų atlyginimų</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>nemokamų atostogų</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>Neoficialiais duomenimis</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>nepadarė</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>nepadarė</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>nepasirodė</i> | V[PrI] | bals-Lietuva-091220-111.txt |
| <i>nepavyko išvengti</i> | VVbendr[prI] | bals-Lietuva-091203-94.txt |
| <i>nerealus ir pareigūnams</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>Nerijui Gudui</i> | Vardai pavardės | bals-Lietuva-091124-85.txt |
| <i>nesirinko</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>neskirs be to</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>nespėjo</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>nesusitaupė</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>nesuspėjo</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>neužgesintos nuorūkos vyras</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>Niekas</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>Niujorko gatvėse</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>Niujorko mero rinkimuose</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>norėjęs</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>nuardė</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>nuimta</i> | Dal[b*][D*](k)D | bals-Lietuva-091213-104.txt |
| <i>numestos</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>Nuogas kaubojus</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>nusikaltimais</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>nusikaltimą</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>nusikaltimo</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>nusiminęs</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>nusprendė pasitraukti iš</i> | VVbendr[prI] | bals-Lietuva-091124-85.txt |
| <i>nusprendęs</i> | Dal[b*][D*](k)D | bals-Pasaulis-090907-2.txt |
| <i>nustebino</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>nutarė</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>nuteistas</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>nutrūko</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>pabandžiusį</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>pagrindo</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>paguldyta</i> | Dal[b*][D*](k)D | bals-Lietuva-091220-111.txt |
| <i>pakeitė</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>palaikantis prokuroras</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>pamatė iš</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>paramos ir labdaros fondo</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>Kalėdų šventėje</i> | | |
| <i>pareigybinės instrukcijas</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>Paryžiaus</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |

| | | |
|---|------------------|-----------------------------|
| <i>parodė</i> | V[Prl] | bals-Lietuva-091124-85.txt |
| <i>pasibaigusios liūdno istorijos pradžia</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>pasirinko</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>pasirodo</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>pasiūlymų</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>paskelbtas Klaipėdos apygardos teismo sprendimas</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>Pastaruojų sprendimu</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>pastatas</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>patikino</i> | V[Prl] | bals-Lietuva-091203-94.txt |
| <i>pavadavo</i> | V[Prl] | bals-Lietuva-091220-111.txt |
| <i>pavargo nuo</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>pažymų</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>pinigų</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>pinigų savivaldybė</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>pirmoji šalies ponis</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>pirmoms gaisrinėms</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Planus nusižudyti</i> | DVBendr | bals-Lietuva-091124-85.txt |
| <i>planų dalyvauti</i> | DVBendr | bals-Pasaulis-090907-2.txt |
| <i>planų įgyvendinti</i> | DVBendr | bals-Lietuva-091124-85.txt |
| <i>policininkams</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>prabangų visureigį</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>pradėjo degti</i> | VVBendr[prl] | bals-Lietuva-091124-85.txt |
| <i>praneša</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>pranešė</i> | V[Prl] | bals-Lietuva-091220-111.txt |
| <i>pranešė</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>pranešimą</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>pravarde</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>prezidento</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>priėmė</i> | V[Prl] | bals-Pasaulis-090907-2.txt |
| <i>priešgaisrinė gelbėjimo tarnyba</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>priešgaisrinės gelbėjimo tarnybos</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>priešgaisrinės priežiūros poskyrio viršininkui</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>privilėgsumokėti už</i> | VVBendr[prl] | bals-Pasaulis-090907-2.txt |
| <i>profesinių sąjungų susivienijimas</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>R. A. Velionės</i> | Vardai pavardės | bals-Lietuva-091124-85.txt |
| <i>R. Matemaitį Švinius</i> | Vardai pavardės | bals-Lietuva-091203-94.txt |
| <i>Raimondą Matemaitį.</i> | Vardai pavardės | bals-Lietuva-091203-94.txt |
| <i>rankos</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>reikalavimus</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>Rinkimams R. Dž.</i> | Vardai pavardės | bals-Pasaulis-090907-2.txt |
| <i>Robertas Džonas Berkas</i> | Vardai pavardės | bals-Pasaulis-090907-2.txt |
| <i>rudenį</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>sako</i> | V[Prl] | bals-Lietuva-091213-104.txt |
| <i>santuokų</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>saugumo</i> | DVBendr | bals-Pasaulis-090907-2.txt |
| <i>sausį</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |

| | | |
|---------------------------------------|------------------|-----------------------------|
| <i>savivaldybės administraciją</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>sekmadienio popietė</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>sekmadienį</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>sekmadienį</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>senutėle</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>sirgusi</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>skirtos baismės</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>skundais</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>sprendimą</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>stogo dangos</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>stovintį medinį namą</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>sugriebė</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>sumažino</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>Suprato</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>surengoje Almos</i> | Dal[b*][D*](k)D | bals-Lietuva-091220-111.txt |
| <i>Susirado</i> | V[PrI] | bals-Lietuva-091124-85.txt |
| <i>sustabdė be</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>sustabdžiusio patrulio</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>sutaupė</i> | V[PrI] | bals-Lietuva-091203-94.txt |
| <i>svarstoma</i> | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| <i>šarvojimo salę</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Šeštadienį Prezidentūroje</i> | Vardai pavardės | bals-Lietuva-091220-111.txt |
| <i>Šiaulių miesto vaikų darželių</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>auklėtojos</i> | | |
| <i>šūkį</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>šventės šeiminkė</i> | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| <i>švietimo</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>Tačiau A.Pukeliui</i> | Vardai pavardės | bals-Lietuva-091203-94.txt |
| <i>tapo</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>Tauragės apskrities viršininko</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>pavaduotoją</i> | | |
| <i>tauragiškio automobilį</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>teigė</i> | V[PrI] | bals-Pasaulis-090907-2.txt |
| <i>teigia</i> | V[PrI] | bals-Lietuva-091213-104.txt |
| <i>teistą</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>tituluojamas nuogu kaubojumi</i> | Dal[b*][D*](k)D | bals-Pasaulis-090907-2.txt |
| <i>trumpikes</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>tūkstantį litų</i> | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| <i>tuometį rajono vicemerą</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>turizmo</i> | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| <i>turįs</i> | Dal[b*][D*](k)D | bals-Pasaulis-090907-2.txt |
| <i>ugniagesiai</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>ūkinius pastatus</i> | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| <i>Užblokavęs kelią</i> | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| <i>V. A. Biržų</i> | Vardai pavardės | bals-Lietuva-091124-85.txt |
| <i>V.Adamkus.</i> | Vardai pavardės | bals-Lietuva-091220-111.txt |
| <i>Vabalninke</i> | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| <i>Vakar Apygardos</i> | Vardai pavardės | bals-Lietuva-091203-94.txt |
| <i>Valdo Adamkaus</i> | Vardai pavardės | bals-Lietuva-091220-111.txt |

| | | |
|--|------------------|-----------------------------|
| Valstybinės priešgaisrinės priežiūros | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| valstybinį kaltinimą | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| vasarį | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| Velionės kūnas | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| vertinamas Tauragės verslininkas | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| Vidutinis ikimokyklinių įstaigų darbuotojų atlyginimas | [sk][b*][d*](k)D | bals-Lietuva-091213-104.txt |
| viešojo saugumo | [sk][b*][d*](k)D | bals-Pasaulis-090907-2.txt |
| viešosios tvarkos pažeidimo | [sk][b*][d*](k)D | bals-Lietuva-091203-94.txt |
| virstančius dūmus | Dal[b*][D*](k)D | bals-Lietuva-091124-85.txt |
| virvė | DVBendr | bals-Lietuva-091124-85.txt |
| virvę | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| žadama | Dal[b*][D*](k)D | bals-Lietuva-091213-104.txt |
| žinomą Arūną | Dal[b*][D*](k)D | bals-Lietuva-091203-94.txt |
| žmona | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |
| žmonai | [sk][b*][d*](k)D | bals-Lietuva-091124-85.txt |
| žmonų | [sk][b*][d*](k)D | bals-Lietuva-091220-111.txt |

9.2. 2 priedas. Eilinio vartotojo rastos frazės

| Phrase | Text |
|--|-----------------------------|
| gaisras | bals-Lietuva-091124-85.txt |
| Kilo | bals-Lietuva-091124-85.txt |
| vyras | bals-Lietuva-091124-85.txt |
| gyvenimo | bals-Lietuva-091124-85.txt |
| Virve | bals-Lietuva-091124-85.txt |
| šarvojimo sale | bals-Lietuva-091124-85.txt |
| Valstybinės priešgaisrinės priežiūros poskyrio viršininkui | bals-Lietuva-091124-85.txt |
| Nerijui Gudui | bals-Lietuva-091124-85.txt |
| norėjo pasikarti | bals-Lietuva-091124-85.txt |
| ženklas | bals-Lietuva-091124-85.txt |
| Klaipėdos apygardos teismo | bals-Lietuva-091203-94.txt |
| Švinius | bals-Lietuva-091203-94.txt |
| sutaupė | bals-Lietuva-091203-94.txt |
| pinigų | bals-Lietuva-091203-94.txt |
| Arūną Pukelį | bals-Lietuva-091203-94.txt |
| kaltinusį | bals-Lietuva-091203-94.txt |
| prokurorą | bals-Lietuva-091203-94.txt |
| viešosios tvarkos | bals-Lietuva-091203-94.txt |
| A.Pukelis | bals-Lietuva-091203-94.txt |
| Tauragės verslininkas | bals-Lietuva-091203-94.txt |
| Tauragės apskrities viršininko pavaduotoją | bals-Lietuva-091203-94.txt |
| Raimondą Matemaitį | bals-Lietuva-091203-94.txt |
| necenzūriniais žodžiais | bals-Lietuva-091203-94.txt |
| 6240 litų baudą | bals-Lietuva-091203-94.txt |
| Auklėtojos | bals-Lietuva-091213-104.txt |
| liko | bals-Lietuva-091213-104.txt |
| atlyginimo | bals-Lietuva-091213-104.txt |

9.3. 3 priedas. Eksperto rastos frazės

| Phrase | |
|--|----------------------------|
| 250 | bals-Pasaulis-090907-2.txt |
| <i>vienos lyties atstovų santuokų</i> | bals-Pasaulis-090907-2.txt |
| <i>dėvėdamas baltas trumpikes</i> | bals-Pasaulis-090907-2.txt |
| <i>biurokratinių reikalavimų</i> | bals-Pasaulis-090907-2.txt |
| <i>dabartiniu Niujorko meru</i> | bals-Pasaulis-090907-2.txt |
| <i>dolerių bauda</i> | bals-Pasaulis-090907-2.txt |
| <i>atsisakymą pateikti pažymą apie jo finansinę padėtį</i> | bals-Pasaulis-090907-2.txt |
| <i>Gatvės muzikantas</i> | bals-Pasaulis-090907-2.txt |
| <i>ir skrybėlę</i> | bals-Pasaulis-090907-2.txt |
| <i>ir visuomeninio transporto naujovių įgyvendinimo</i> | bals-Pasaulis-090907-2.txt |
| <i>kandidatūrą</i> | bals-Pasaulis-090907-2.txt |
| <i>kaubojiškus batus</i> | bals-Pasaulis-090907-2.txt |
| <i>liepos mėnesį</i> | bals-Pasaulis-090907-2.txt |
| <i>mažesnėmis galimybėmis</i> | bals-Pasaulis-090907-2.txt |
| <i>mero postą</i> | bals-Pasaulis-090907-2.txt |
| <i>miesto vadovo rinkimuose</i> | bals-Pasaulis-090907-2.txt |
| <i>pasiūlymų dėl mokesčių,</i> | bals-Pasaulis-090907-2.txt |
| <i>naujienu agentūra</i> | bals-Pasaulis-090907-2.txt |
| <i>Niekas</i> | bals-Pasaulis-090907-2.txt |
| <i>Niujorko gatvėse</i> | bals-Pasaulis-090907-2.txt |
| <i>Niujorko mero rinkimuose</i> | bals-Pasaulis-090907-2.txt |
| <i>Nuogas kaubojus</i> | bals-Pasaulis-090907-2.txt |
| <i>pažymą</i> | bals-Pasaulis-090907-2.txt |
| <i>sekmadienį</i> | bals-Pasaulis-090907-2.txt |
| <i>sprendimą</i> | bals-Pasaulis-090907-2.txt |
| <i>šūkį</i> | bals-Pasaulis-090907-2.txt |
| <i>turizmo</i> | bals-Pasaulis-090907-2.txt |
| <i>viešojo saugumo</i> | bals-Pasaulis-090907-2.txt |
| <i>būtina įvykdyti</i> | bals-Pasaulis-090907-2.txt |
| <i>Buvusio kandidato teigimu</i> | bals-Pasaulis-090907-2.txt |
| <i>nusprendęs kandidatuoti mero rinkimuose</i> | bals-Pasaulis-090907-2.txt |
| <i>tituluojamas nuogu kaubojumi</i> | bals-Pasaulis-090907-2.txt |
| <i>turįs daug pasiūlymų dėl mokesčių</i> | bals-Pasaulis-090907-2.txt |
| <i>atsisakymą pateikti pažymą apie jo finansinę padėtį</i> | bals-Pasaulis-090907-2.txt |
| <i>planų dalyvauti</i> | bals-Pasaulis-090907-2.txt |
| <i>JAV</i> | bals-Pasaulis-090907-2.txt |
| <i>atsisakė</i> | bals-Pasaulis-090907-2.txt |
| <i>yra</i> | bals-Pasaulis-090907-2.txt |
| <i>Paskutiniu lašu</i> | bals-Pasaulis-090907-2.txt |
| <i>nepadarė</i> | bals-Pasaulis-090907-2.txt |
| <i>nutarė nedalyvaut</i> | bals-Pasaulis-090907-2.txt |
| <i>pasirinko</i> | bals-Pasaulis-090907-2.txt |
| <i>pasirodo</i> | bals-Pasaulis-090907-2.txt |
| <i>pavargo nuo</i> | bals-Pasaulis-090907-2.txt |
| <i>praneša</i> | bals-Pasaulis-090907-2.txt |
| <i>pranešė</i> | bals-Pasaulis-090907-2.txt |

| | |
|---|----------------------------|
| <i>priėmė</i> | bals-Pasaulis-090907-2.txt |
| <i>tapo</i> | bals-Pasaulis-090907-2.txt |
| <i>teigė</i> | bals-Pasaulis-090907-2.txt |
| <i>AP.</i> | bals-Pasaulis-090907-2.txt |
| <i>Maiklu Blumbergu</i> | bals-Pasaulis-090907-2.txt |
| <i>Rinkimams</i> | bals-Pasaulis-090907-2.txt |
| <i>Robertas Džonas Berkas</i> | bals-Pasaulis-090907-2.txt |
| <i>ketina konkuruoti su</i> | bals-Pasaulis-090907-2.txt |
| <i>privalės sumokėti už</i> | bals-Pasaulis-090907-2.txt |
| <i>R.Dž.Berkas</i> | bals-Pasaulis-090907-2.txt |
| <i>R.Dž.Berko</i> | bals-Pasaulis-090907-2.txt |
| <i>16 kvadratinių metrų kambarys</i> | bals-Lietuva-091124-85.txt |
| <i>17.13 val.</i> | bals-Lietuva-091124-85.txt |
| <i>58 metų vyras V. A.</i> | bals-Lietuva-091124-85.txt |
| <i>6 kvadratinius metrus stogo dangos</i> | bals-Lietuva-091124-85.txt |
| <i>Aiškino</i> | bals-Lietuva-091124-85.txt |
| <i>atvira liepsna</i> | bals-Lietuva-091124-85.txt |
| <i>Biržų priešgaisrinė gelbėjimo tarnyba</i> | bals-Lietuva-091124-85.txt |
| <i>Biržų priešgaisrinės gelbėjimo tarnybos Valstybinės</i> | bals-Lietuva-091124-85.txt |
| <i>priešgaisrinės priežiūros poskyrio viršininkui Nerijui</i> | |
| <i>Gudui</i> | |
| <i>buvę baldai</i> | bals-Lietuva-091124-85.txt |
| <i>buvo išvežtas į</i> | bals-Lietuva-091124-85.txt |
| <i>buvo pasiruošęs.</i> | bals-Lietuva-091124-85.txt |
| <i>Buvo taip nusiminęs</i> | bals-Lietuva-091124-85.txt |
| <i>degė</i> | bals-Lietuva-091124-85.txt |
| <i>per du metrus nuo degusio namo stovintį medinį</i> | bals-Lietuva-091124-85.txt |
| <i>namą.</i> | |
| <i>Gaisras</i> | bals-Lietuva-091124-85.txt |
| <i>Gaisrininkai</i> | bals-Lietuva-091124-85.txt |
| <i>Gaisro metu</i> | bals-Lietuva-091124-85.txt |
| <i>gavo</i> | bals-Lietuva-091124-85.txt |
| <i>gyvenimo</i> | bals-Lietuva-091124-85.txt |
| <i>iš kambario virstančius dūmus</i> | bals-Lietuva-091124-85.txt |
| <i>išdegė</i> | bals-Lietuva-091124-85.txt |
| <i>Išgelbėjo</i> | bals-Lietuva-091124-85.txt |
| <i>įgyvendinti nespėjo</i> | bals-Lietuva-091124-85.txt |
| <i>įgyvendinti nesuspėjo</i> | bals-Lietuva-091124-85.txt |
| <i>kambaryje</i> | bals-Lietuva-091124-85.txt |
| <i>Ketinių</i> | bals-Lietuva-091124-85.txt |
| <i>kilęs gaisras</i> | bals-Lietuva-091124-85.txt |
| <i>Kilo</i> | bals-Lietuva-091124-85.txt |
| <i>kilo nuo</i> | bals-Lietuva-091124-85.txt |
| <i>liepsnas</i> | bals-Lietuva-091124-85.txt |
| <i>likimo atsiųstas ženklas</i> | bals-Lietuva-091124-85.txt |
| <i>Manoma</i> | bals-Lietuva-091124-85.txt |
| <i>mirė</i> | bals-Lietuva-091124-85.txt |
| <i>namas</i> | bals-Lietuva-091124-85.txt |
| <i>nebuvo</i> | bals-Lietuva-091124-85.txt |

| | |
|---|-----------------------------|
| <i>dirbs be</i> | bals-Lietuva-091213-104.txt |
| <i>Esq</i> | bals-Lietuva-091213-104.txt |
| <i>Gavo</i> | bals-Lietuva-091213-104.txt |
| <i>gruodj</i> | bals-Lietuva-091213-104.txt |
| <i>išnaudotas metų mokus fondas</i> | bals-Lietuva-091213-104.txt |
| <i>jstaigų vadovėms</i> | bals-Lietuva-091213-104.txt |
| <i>jteikė</i> | bals-Lietuva-091213-104.txt |
| <i>Yra</i> | bals-Lietuva-091213-104.txt |
| <i>yra apie</i> | bals-Lietuva-091213-104.txt |
| <i>kreipėsi j</i> | bals-Lietuva-091213-104.txt |
| <i>Lapkritj</i> | bals-Lietuva-091213-104.txt |
| <i>liko be</i> | bals-Lietuva-091213-104.txt |
| <i>Metu</i> | bals-Lietuva-091213-104.txt |
| <i>metų viduryje</i> | bals-Lietuva-091213-104.txt |
| <i>miesto biudžete</i> | bals-Lietuva-091213-104.txt |
| <i>Miesto darželių vedėjos</i> | bals-Lietuva-091213-104.txt |
| <i>mokos fondus</i> | bals-Lietuva-091213-104.txt |
| <i>nebeliko</i> | bals-Lietuva-091213-104.txt |
| <i>Nebesitiki</i> | bals-Lietuva-091213-104.txt |
| <i>nemokamų atlyginimų</i> | bals-Lietuva-091213-104.txt |
| <i>nemokamų atostogų</i> | bals-Lietuva-091213-104.txt |
| <i>nesirinko</i> | bals-Lietuva-091213-104.txt |
| <i>neskirs</i> | bals-Lietuva-091213-104.txt |
| <i>nesusitaupė</i> | bals-Lietuva-091213-104.txt |
| <i>nuimta po</i> | bals-Lietuva-091213-104.txt |
| <i>pinigų</i> | bals-Lietuva-091213-104.txt |
| <i>reikalavimus</i> | bals-Lietuva-091213-104.txt |
| <i>Rudenj</i> | bals-Lietuva-091213-104.txt |
| <i>Sako</i> | bals-Lietuva-091213-104.txt |
| <i>sausj</i> | bals-Lietuva-091213-104.txt |
| <i>Savivaldybė</i> | bals-Lietuva-091213-104.txt |
| <i>savivaldybės administraciją</i> | bals-Lietuva-091213-104.txt |
| <i>sumažino</i> | bals-Lietuva-091213-104.txt |
| <i>Šiaulių ikimokyklinių įstaigų darbuotojos</i> | bals-Lietuva-091213-104.txt |
| <i>Šiaulių miesto vaikų darželių auklėtojos</i> | bals-Lietuva-091213-104.txt |
| <i>Šiaulių švietimo darbuotojų profesinių sąjungų</i> | bals-Lietuva-091213-104.txt |
| <i>susivienijimas</i> | bals-Lietuva-091213-104.txt |
| <i>Teigia</i> | bals-Lietuva-091213-104.txt |
| <i>tūkstantį litų</i> | bals-Lietuva-091213-104.txt |
| <i>vasarj</i> | bals-Lietuva-091213-104.txt |
| <i>Vidutinis ikimokyklinių įstaigų darbuotojų</i> | bals-Lietuva-091213-104.txt |
| <i>atlyginimas</i> | bals-Lietuva-091213-104.txt |
| <i>vieną atlyginimą</i> | bals-Lietuva-091213-104.txt |
| <i>žadama išmokėti</i> | bals-Lietuva-091213-104.txt |