

# **Lietuvių kalbos teksto sintaksinės-semantinės analizės informacinė sistema**

## **LKSSAIS vystymas**

Modernizuojamos LKSSAIS Lietuviškų dokumentų tekstų pagrindinės  
ir specializuotos statistinės analizės IT sprendimo komponento  
administratoriaus instrukcija

**Lietuvių kalbos teksto sintaksinės-semantinės analizės informacinė sistema**

## TURINYS

<b>1. Įžanga .....</b>	<b>4</b>
1.1. Dokumento paskirtis.....	4
1.2. Sutrumpinimai .....	4
1.3. Panaudotų dokumentų sąrašas.....	4
1.4. Statistikos IT sprendimo paskirtis ir tikslai.....	4
<b>2. Statistikos IT sprendimo administratoriaus instrukcija .....</b>	<b>5</b>
2.1. Statistikos IT sprendimo komponento sandara .....	5
2.2. Statistikos IT sprendimo komponento diegimas į sistemą.....	5
2.2.1. Docker konteinerio sutvėrimas .....	5
2.2.2. Diegimas į Kubernetes aplinką.....	5
2.3. Statistikos IT sprendimo komponento šalinimas .....	6
2.3.1. Docker konteinerio pašalinimas .....	6
2.3.2. Pašalinimas iš Kubernetes aplinkos.....	6
2.4. Statistikos IT sprendimo komponento veikimas .....	7
2.4.1. Įėjimas .....	7
2.4.2. Išėjimas .....	7
<b>3. Priedai .....</b>	<b>8</b>
3.1. Modernizuojamos LKSSAIS Statistiko IT sprendimo komponento skaičiuojami statistikos parametrai (atkelta iš [2]) .....	8
<b>4. Dokumento istorija .....</b>	<b>Error! Bookmark not defined.</b>

## LENTELIŲ SĄRAŠAS

1.1 lentelė. Sutrumpinimai .....	4
1.2 lentelė. Panaudotų dokumentų sąrašas .....	4
2.1 lentelė. Komponento tinklinės paslaugos metodai .....	5
3.1 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo komponento skaičiuojami statistikos parametrai (atkelta iš [2]) .....	8
3.2 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo komponento skaičiuojami morfologinių savybių statistikos parametrai (atkelta iš [2]) .....	9
3.3 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo skaičiuojami skaitomumo indeksai ir susijusios statistikos (atkelta iš [2], papildyta komentarais).....	10
4.1 lentelė. Dokumentų keitimo istorija .....	<b>Error! Bookmark not defined.</b>

# 1. Įžanga

## 1.1. Dokumento paskirtis

Šio dokumento paskirtis – pateikti modernizuojamos LKSSAIS Statistikos IT sprendimo komponento administravimo vadovą.

## 1.2. Sutrumpinimai

1.1 lentelė. Sutrumpinimai

Santrumpa	Paiškinimas
LKSSAIS	Lietuvių kalbos sintaksinės ir semantinės analizės informacinė sistema
IT	Informacinės technologijos

## 1.3. Panaudotų dokumentų sąrašas

1.2 lentelė. Panaudotų dokumentų sąrašas

Eil. Nr.	Pavadinimas
1.	Modernizuojamos LKSSAIS Statistikos IT sprendimo komponento detalios analizės ir projektavimo specifikacija
2.	TEKSTO PILNOS IR SPECIALIZUOTOS STATISTINĖS ANALIZĖS BENDROJO NAUDOJIMO IT SPRENDIMAS, Konceptijos specifikacija, Kaunas, 2018.

## 1.4. Statistikos IT sprendimo paskirtis ir tikslai

Lietuviškų dokumentų tekstų pagrindinės ir specializuotos statistinės analizės IT sprendimas (toliau Statistikos IT sprendimas) analizuos dokumento teksto leksinius segmentus ir morfologines leksinių vienetų žymas ir suskaičiuos bazinius teksto statistikos parametrus (angl. *basic text statistics*), dalinę kalbos dalių statistiką, teksto leksinį tankį (angl. *lexical density*). Suskaičiuoti parametrai bus grąžinami IT sprendimo funkcionalumą iškvietusiai sistemai, posistemei ar komponentui.

## 2. Statistikos IT sprendimo administratoriaus instrukcija

### 2.1. Statistikos IT sprendimo komponento sandara

Statistikos IT sprendimo komponentas veikia Docker konteineryje.

Komponento vykdomieji failai konteineryje patalpinti kataloge */trm/webservice*.

Komponento tinklinė paslauga konteineryje pasiekama adresu *http://0.0.0.0:4000/api*. Tinklinės paslaugos kontrakto aprašymą *Swagger (OpenAPI)* formatu galima gauti kreipiantis adresu *http://0.0.0.0:4000/swagger/api/swagger.json*.

Komponento tinklinės paslaugos adresas įdiegus Kubernetes aplinkoje priklauso nuo to kokių adresu Kubernetes išviešinta konteineryje veikianti tinklinė paslauga. Šio adreso parinkimas paliekamas administratoriui. Administratorius gali patikrinti ar parinktas adresas veikia teisingai per jį kviesdamas vidinį konteinerio resursą *http://0.0.0.0:4000/api/ping*. Jeigu resursas atsako formatu „*PONG: laiko-štampos*“, vidinė paslauga pasiekama.

Komponento tinklinė paslauga turi du metodus, aprašytus lentelėje 2.1.

2.1 lentelė. Komponento tinklinės paslaugos metodai

Metodas	Paaiškinimas
<i>/api/</i>	Skaičiuoja duoto dokumento statistinius ir skaitomumo įverčius. Kviečiamas per HTTP POST antraštėje <i>Content-Type</i> nurodant turinio tipą <i>application/json</i> bei teisingą perduodamo teksto koduotę. Įėjimo ir išėjimo formatai aprašyti skyriuje 2.3
<i>/api/ping</i>	Leidžia patikrinti ar paslauga veikia. Kviečiamas per HTTP GET arba HTTP POST. Atsako formatu „ <i>PONG: laiko-štampos</i> “.

### 2.2. Statistikos IT sprendimo komponento diegimas į sistemą

#### 2.2.1. Docker konteinerio sutvėrimas

Komponento Docker konteinerio sutvėrimui reikia turėti komponento išeitinius tekstus bei vykdomąją aplinką tenkinančia tokius parametrus: *Ubuntu Linux v>=18.04*, *.NET Core v>=2.1*, *Docker v>=18*.

Komponento išeitiniuose tekstuose reikia eiti į katalogą *TextReadability/WebService* ir jame įvykdyti komandas:

```
dotnet publish
docker build -t "trm:0.1" ./
```

Lokalioje sistemoje bus sutvėrta komponento Docker saugykla (angl. *repository*) pavadinta *trm* su žyma (angl. *tag*) *0.1*.

#### 2.2.2. Diegimas į Kubernetes aplinką

Komponento diegimui į Kubernetes aplinką reikia turėti komponento Docker konteinerį patalpintą Kubernetes aplinkai pasiekiamoje Docker saugykloje (angl. *repository*).

Diegiant Kubernetes aplinkoje komponentui sukuriama diegimo (angl. *deployment*) bei paslaugos (angl. *service*) aprašai. Diegimo aprašo pavyzdys:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: semantika2-trm-webservice
  labels:
    app: semantika2-trm-webservice
spec:
  replicas: 1
  selector:
```

```

matchLabels:
  app: semantika2-trm-webservice
template:
  metadata:
    labels:
      app: semantika2-trm-webservice
  spec:
    containers:
      - name: semantika2-trm-webservice
        image: docker-konteinerio-adresas
        ports:
          - containerPort: 4000

```

čia *docker-konteinerio-adresas* yra komponento Docker konteinerio adresas Kubernetes aplinkai pasiekiamoje Docker saugykloje.

Laikant, kad diegimo aprašas patalpintas faile *diegimas.yml*, Kubernetes aplinkoje diegimas sukuriamas naudojant komandą:

```
kubectl apply -f diegimas.yml
```

Paslaugos aprašo pavyzdys:

```

apiVersion: v1
kind: Service
metadata:
  name: semantika2-trm-webservice
spec:
  ports:
    - port: 80
      protocol: TCP
      targetPort: 4000
  selector:
    app: semantika2-trm-webservice

```

šis aprašas Kubernetes aplinkoje komponento tinklinę paslaugą padarys prieinamą per TCP prievadą 80. Prievadą galima pakeisti į bet kurį norimą.

Laikant, kad paslaugos aprašas patalpintas faile *paslauga.yml*, Kubernetes aplinkoje paslauga sukuriamą naudojant komandą:

```
kubectl apply -f paslauga.yml
```

## 2.3. Statistikos IT sprendimo komponento šalinimas

### 2.3.1. Docker konteinerio pašalinimas

Šiame skyriuje laikoma, kad komponento Docker konteineris lokaliaje aplinkoje buvo sutvertas laikantis skyriuje 2.2.1 nurodytų žingsnių.

Iš lokalsios aplinkos komponento Docker konteineris pašalinimas naudojant komandą:

```
docker image remove trm:0.1
```

### 2.3.2. Pašalinimas iš Kubernetes aplinkos

Šiame skyriuje laikoma, kad komponentas Kubernetes aplinkoje buvo įdiegtas laikantis skyriuje 2.2.2 nurodytų žingsnių.

Komponento paslauga pašalinama naudojant komandą:

```
kubectl delete service semantika2-trm-webservice
```

Komponento diegimas pašalinimas naudojant komandą:

```
kubectl delete deployment semantika2-trm-webservice
```

## 2.4. Statistikos IT sprendimo komponento veikimas

### 2.4.1. Įėjimas

Komponento įėjimas turi būti JSON dokumentas tokiu formatu:

```
{
  "body": "dokumento tekstas",
  "annotations": {
    "lex" : ... leksikografinės anotacijos dokumentui ...,
    "morph" : ... morfologinės anotacijos dokumentui ...
  }
}
```

Bet kokie papildomi laukai įėjimo dokumente bus ignoruojami.

### 2.4.2. Išėjimas

Klaidos atveju komponentas grąžina HTTP atsakymo kodą 500 bei išėjimą tokiu formatu:

```
{
  "error": {
    "status" : 500,
    "message" : "klaidos žinutė"
  }
}
```

Komponento išėjimas normaliu atveju yra JSON dokumentas tokiu formatu:

```
{
  "text": { parametrai 1-7, 3.1 lentelėje },
  "lexical": { parametrai 8-13, 3.1 lentelėje },
  "means": { parametrai 14-20, 3.1 lentelėje },
  "min": { parametrai 21-23, 3.1 lentelėje },
  "max": { parametrai 24-26, 3.1 lentelėje },
  "word": { parametrai 27-28, 3.1 lentelėje },
  "freq": { parametrai 29-32, 3.1 lentelėje },
  "morph": { parametrai 33-44, 3.1 lentelėje },
  "morphParameters": { visi parametrai 3.2 lentelėje },
  "readability": { skaitomumo indeksų įverčiai 3.3 lentelėje }
}
```

Čia kiekvienas parametras ir skaitomumo indeksas yra atskiras JSON struktūros laukas, kurio pavadinimas atspindi atitinkamo parametro arba skaitomumo indekso pavadinimą. Detalią dokumento struktūrą galima gauti komponento tinklinės paslaugos *Swagger (OpenAPI)* dokumentacijoje kreipiantis į resursą `/swagger/api/swagger.json` , arba tiesiog į `/swagger/index.html` .

### 3. Priedai

#### 3.1. Modernizuojamos LKSSAIS Statistiko IT sprendimo komponento skaičiuojami statistikos parametrai (atkelta iš [2])

3.1 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo komponento skaičiuojami statistikos parametrai (atkelta iš [2])

Nr.	Aprašymas
	<b>Teksto statistika</b>
1.	Simbolių skaičius su tarpais
2.	Simbolių skaičius be tarpų
3.	Skiemenų skaičius
4.	Leksemų skaičius
5.	Žodžių skaičius
6.	Unikalių žodžių (lemų) skaičius
7.	Sakinių skaičius
	<b>Leksinis tankis</b>
8.	Viso teksto Leksinių (turinį perteikiančių) žodžių skaičius
9.	Viso teksto Neleksinių (funkcinių) žodžių skaičius
10.	Viso teksto Leksinis tankis
11.	Sakinio leksinių (turinį perteikiančių) žodžių skaičius
12.	Sakinio neleksinių (funkcinių) žodžių skaičius
13.	Sakinio leksinis tankis
	<b>Vidurkiai</b>
14.	Žodžio ilgio vidurkis simboliais
15.	Žodžio ilgio vidurkis skiemenimis
16.	Žodžio ilgio mediana (simboliais/skiemenimis)
17.	Standartinis žodžio ilgio nuokrypis
18.	Sakinio ilgio vidurkis žodžiais (leksemomis)
19.	Sakinio ilgio mediana
20.	Sakinio ilgio standartinis nuokrypis
	<b>Min</b>
21.	Trumpiausias žodis simboliais
22.	Trumpiausias žodis skiemenimis
23.	Trumpiausias sakiny s žodžiais (leksemomis)
	<b>Max</b>
24.	Ilgiausias žodis simboliais
25.	Ilgiausias žodis skiemenimis
26.	Ilgiausias sakiny s žodžiais (leksemomis)
	<b>Žodžių ilgio simboliais statistika</b>
27.	Žodžių skaičius ir procentinė dalis kiekvienam žodžio ilgiui simboliais (t. y., vieno simbolio žodžių skaičius ir procentinė dalis, dviejų simbolių žodžių skaičius ir procentinė dalis, trijų simbolių žodžių skaičius ir procentinė dalis ir t.t. iki didžiausio, koks yra analizuotame tekste, ilgio simboliais žodžių skaičiaus ir procentinės dalies)
	<b>Žodžių ilgio skiemenimis statistika</b>
28.	Žodžių skaičius ir procentinė dalis kiekvienam žodžio ilgiui skiemenimis (t. y., vieno skiemens žodžių skaičius ir procentinė dalis, dviejų skiemenų žodžių skaičius ir procentinė dalis, trijų skiemenų žodžių skaičius ir procentinė dalis ir t.t. iki didžiausio, koks yra analizuotame tekste, ilgio skiemenimis žodžių skaičiaus ir procentinės dalies)
	<b>Dažniai</b>
29.	Dažniausių žodžių (funkcinių ir nefunkcinių)
30.	Dažniausių nefunkcinių žodžių pagal dydį ir spalvą (debesis)
31.	Dažniausių žodžių porų
32.	Dažniausių žodžių tripletų
	<b>Teksto morfologinė struktūra</b>



Nr.	Aprašymas
33.	Daiktavardžių skaičius ir procentinė dalis
34.	Asmenuojamų veiksmažodžių skaičius ir procentinė dalis
35.	Bendratis veiksmažodžių skaičius ir procentinė dalis
36.	Dalyvių skaičius ir procentinė dalis
37.	Padalyvių skaičius ir procentinė dalis
38.	Pusdalyvių skaičius ir procentinė dalis
39.	Būdvardžių skaičius ir procentinė dalis
40.	Skaitvardžių skaičius ir procentinė dalis
41.	Įvardžių skaičius ir procentinė dalis
42.	Prieveiksmių skaičius ir procentinė dalis
43.	Kitų kalbos dalių skaičius ir procentinė dalis
44.	Neatpažintų žodžių leksemų skaičius ir procentinė dalis

3.2 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo komponento skaičiuojami morfologinių savybių statistikos parametrai (atkelta iš [2])

Kalbos dalis	Vertinama morfologinė savybė	
	pavadinimas	galimos reikšmės
daiktavardis	pobūdis	bendrinis, tikrinis
	giminė	bendroji, moteriška, vyriška
	skaičius	daugiskaita, vienaskaita, dviskaita
	linksnis	vardininkas, kilmininkas, naudininkas, galininkas, įnagininkas, vietininkas, šauksmininkas
	ar sangražinis	sangražinis, nesangražinis
	vardo eilė	vardas, pavardė, vietovė
asmenuojamas veiksmažodis	laikas	esamasis, būtasis kartinis, būtasis dažninis, būsimasis, -
	asmuo	I, II, III
	skaičių	daugiskaita, vienaskaita
	ar neigiamas	teigiamas, neigiamas
	ar sangražinis	sangražinis, nesangražinis
	nuosaka	tiesioginė, tariamoji, liepiamoji
bendratis veiksmažodžiams	ar neigiamas	teigiamas, neigiamas
	ar sangražinis	sangražinis, nesangražinis
dalyviai	laikas	esamasis, būtasis kartinis, būtasis dažninis, būsimasis, -
	giminė	bevardė, moteriška, vyriška
	skaičius	daugiskaita, vienaskaita, -
	rūšis	veikiamoji, neveikiamoji, reikiamoji
	ar neigiamas	teigiamas, neigiamas
	ar įvardžiuotinis	įvardžiuotinis, neįvardžiuotinis
	linksnis	vardininkas, kilmininkas, naudininkas, galininkas, įnagininkas, vietininkas, šauksmininkas
	ar sangražinis	sangražinis, nesangražinis
padalyviai	laikas	esamasis, būtasis kartinis, būtasis dažninis, būsimasis
	ar neigiamas	teigiamas, neigiamas
	ar sangražinis	sangražinis, nesangražinis
pusdalyviai	skaičius	daugiskaita, vienaskaita
	giminė	moteriška, vyriška
	ar neigiamas	teigiamas, neigiamas
	ar sangražinis	sangražinis, nesangražinis
būdvardžiai	laipsnis	nelyginis, aukštesnis, aukščiausias, -
	giminė	bevardė, moteriška, vyriška, -
	skaičius	daugiskaita, vienaskaita, dviskaita, -

Kalbos dalis	Vertinama morfologinė savybė	
	pavadinimas	galimos reikšmės
skaitvardžiai	linksnis	vardininkas, kilmininkas, naudininkas, galininkas, įnagininkas, vietininkas, šauksmininkas
	ar įvardžiuotinis	įvardžiuotinis, neįvardžiuotinis, -
	grupė	kiekiniai, kuopiniai, -
	giminė	bevardė, moteriška, vyriška, -
	skaičius	daugiskaita, vienaskaita, -
	linksnis	vardininkas, kilmininkas, naudininkas, galininkas, įnagininkas, vietininkas, šauksmininkas, -
	forma	raidinė, skaitmeninė, romėniškas
įvardžiai	ar įvardžiuotinis	įvardžiuotinis, neįvardžiuotinis, -
	giminė	bevardė, moteriška, vyriška, -
	skaičius	daugiskaita, vienaskaita, dviskaita, -
	linksnis	vardininkas, kilmininkas, naudininkas, galininkas, įnagininkas, vietininkas, šauksmininkas
prieveiksmiai	ar įvardžiuotinis	įvardžiuotinis, neįvardžiuotinis, -
	laipsnis	nelyginis, aukštesnis, aukščiausias
kitos kalbos dalys	pobūdis, jei tokį turi	jaustukas, išiktukas, dalelytė, prielinksnis, jungtukas, trumpinys, skyrybos ženklas, neatpažintas

3.3 lentelė. Modernizuojamos LKSSAIS Statistikos IT sprendimo skaičiuojami skaitomumo indeksai ir susijusios statistikos (atkelta iš [2], papildyta komentarais)

Nr.	Aprašymas
	<b>Skaitomumo indeksai</b>
1.	<i>Flesch-Kincaid</i> indeksas (skyla į du indeksus: <i>Flesch-Kincaid Grade Level</i> ir <i>Flesch-Kincaid Reading Ease</i> )
2.	<i>Gunning Fog</i> indeksas
3.	<i>Coleman-Liau</i> indeksas
4.	<i>SMOG</i> indeksas
5.	<i>Automated Readability</i> indeksas
6.	<i>Indeksų vidurkis</i> (apima 1-5 eilutes, 8-9 yra metrikos)
7.	<i>Indeksų mediana</i> (apima 1-5 eilutes, 8-9 yra metrikos)
8.	<i>Fry</i> skaitomumo metrika (X ir Y dedamosios <i>Fry</i> grafike)
9.	<i>Raygor</i> skaitomumo metrika (X ir Y dedamosios <i>Raygor</i> grafike)