

Marķētājs ir pusautomātisks programmrīks, kas paredzēts latviešu valodas tekstu morfoloģiskai un sintaktiskai analīzei. Ar tā palīdzību katrai vārdformai tiek dots morfoloģiskais raksturojums un, ja iespējams, noteikta arī vārda sintaktiskā loma teikumā.

Morfoloģiskajam raksturojumam tiek izmantoti speciāli simboli. Marķētājā tiek izmantota latviešu valodas morfoloģisko pazīmju kopa, kas ir atvasināta no *MULTEXT-East* pazīmju kopas¹. Morfoloģisko pazīmju kopu var aplūkot arī, ejot uz izvēlni [Palīgs] \ [Morfoloģisko pazīmju kopa].

Marķējamā teksta atvēršana

Marķējamo tekstu var atvērt, ejot uz izvēlni [Korpuss] \ [Atvērt marķējamo tekstu]

Marķētā teksta saglabāšana

Marķēto tekstu var saglabāt, ejot uz izvēlni [Korpuss] \ [Saglabāt marķēto tekstu]. Marķētie dati tiek saglabāti gan teksta formātā, gan .xml formātā. Saglabātajā teksta failā vārdi sakārtoti viens zem otra, un blakus tiem dots marķējums.

Teksta fragmentēšana

Atverot marķējamo tekstu, tas automātiski tiek sadalīts fragmentos un vārdformām tiek pievienots morfoloģiskais raksturojums/morfoloģiskās pazīmes, t. i., tiek piedāvāti visi iespējamie marķējuma varianti. Teksta fragmenti ir sakārtoti viens zem otra saskarnes kreisajā pusē (sk. 1. att.).

Teksta fragmentu krāsa. Dzeltēnā krāsā iekrāsoti tie teksta fragmenti, kuru marķēšana ir iesākta, bet nav pabeigta. Savukārt pelēkā krāsā ir tie teksta fragmenti, kuru marķēšana ir pabeigta. Teksta fragments, kas konkrētajā brīdī tiek analizēts un kura vārdformas un pieturzīmes ir dotas marķējuma tabulā, ir zaļā krāsā (sk. 1. att.).

Teksta fragmentus iespējams **dzēst** [Dzēst rindiņu] un **apvienot** [Apvienot ar nākamo rindiņu], starp teksta fragmentiem var **ievietot tukšu rindiņu** [Iesprausts tukšu rindiņu]. To var izdarīt, novietojot kursoru uz teksta fragmenta un uzklikšķinot uz tā ar peles labo taustiņu. Parādās izvēlnes logs (sk. 2. att.), kurā

¹ Sk. sadaļu [Rīki] \ [Morfoloģiskais analizators] adresē <http://www.semti-kamols.lv>

jāizvēlas opcija [Apvienot ar nākamo rindiņu], [Dzēst rindiņu] vai [Iespraust tukšu rindiņu].

Teksta fragmenti

Vārdformas un to morfosintaktiskās pazīmes

Vārdformas *ūdēns* morfoloģiskās analīzes varianti

1. att. Markētāja saskarne

Ielikt tekstu no starpliktuves

Turpināt markēt no šīs rindiņas

Pārmarkēt šo rindiņu

Apvienot ar nākamo rindiņu

Iespraust tukšu rindiņu

Dzēst rindiņu

Izvēlēties citu analīzes variantu

2. att. Teksta fragmentu apstrādes komandriņas

Morfoloģiskā analīze

Atverot markējamo tekstu, visas pirmā teksta fragmenta vārdformas un pieturzīmes tiek sarakstītas *Markējuma tabulā* saskarnes labajā pusē (sk. 3. att. A) un tām tiek piešķirts kārtas numurs (sk. 3. att. B). Automātiski tiek piedāvāti pirmās vārdformas vai pieturzīmes analīzes varianti (sk. 3. att. C)

Markējuma tabulas (sk. 3. att.) 1. ailē doti vārdformu un pieturzīmju kārtas numuri, 2. ailē *Vārds* — vārdformas un pieturzīmes, savukārt 3. ailē *Markējums* paredzēta markējuma / morfoloģisko pazīmju virknei, 4. ailē *Pamatforma* — vārda pamatformai, 5. ailē *Loma* paredzēta sintaktiskās lomas atspoguļošanai, 6. ailē *Paskaidro* norāda uz vārdu, kam attiecīgā vārdforma pakārtota.

Markētājs - 0.9.59 [C:\Users\Iļze Auziņa\Desktop\annotetajs 22.05.09\annotator-r592\demo1.txt]

Korpuss Rīki Palīgs

Čānks

Markējuma tabula Sintakse

Dabiski, ka ūdens nevar pārvērsties par zivi vai taureni.

Ūdens vispār nevar pārvērsties.

Tīrs ūdens mūžīgi ir tīrs ūdens.

Tātad Parmenīdam bija taisnība, ka "nekas nepārvēršas".

Tajā pašā laikā Empedokls bija vienprātis ar Heraklītu, ka mums jāpajaujas.

Mums jātic tam, ko redzam, un mēs redzam tieši to, ka.

dabā viss pastāvīgi mainās.

Empedokls secināja, ka jāatmet pieņēmums par vienas pirmvielas pastāvēšanu.

Ne ūdens, ne gaiss paši par sevi nevar pārvērsties par rožu krūmu vai taureni.

Tātad nav iespējams, ka dabā ir tikai viena viela, kas ir visa pamatā.

Empedokls apgalvoja, ka dabā pastāv pavisam četri pamatelementi jeb "cēloņi".

Tie, viņaprāt, bija zeme, gaiss, uguns un ūdens.

Visas pārvērtības dabā nāca no tā, ka šie četri elementi savienojās un atkal sadalījās.

Jo viss sastāv no zemes, gaisa, uguns un ūdens, tikai dažādās attiecībās.

Kad zieds novīst vai kāds

✓ Rādīt visus vārdus

Npk.	Vārds	Markējums	Pamatforma	Loma	Paskaidro
1	Dabiski				
2	,				
3	ka				
4	ūdens				
5	nevar				
6	pārvērsties				
7	par				
8	zivi				
9	vai				
10	taureni				
11	.				

Vārds	Markējums	Galotnes...	Ticamība	Apraksts
+ Dabiski	<[a,f,m,p,n,n,p], 'dabiski'>	140	1	Ipašības vārds, 140, Kādības, dabiski, 3, Nominatīvs, dabiski, 12, 18627, Daudzskaitlis, N...
+ Dabiski	<[r, _], 'dabiski'>	179	1	Kādības, 179, 12, Apstākļa vārds, dabiski, 18627, 0, dabiski, Valērija leksikons, r...
+ Dabiski	Es zinu labāk!			
+ Dabiski	Reziduālis			

3. att. Markētāja saskarne: analizējamās vārdformas un vārdformas *dabiski* analīzes varianti

Katru markējamo vārdformu vai pieturzīmi var aplūkot tabulā saskarnes lejas daļā: tur redzama vārdforma (vai pieturzīme), iespējamais morfoloģiskais markējums un markējuma apraksts, norādīta ticamība un galotnes numurs (sk. 3. att. C). Izvēlēties atbilstošo markējumu iespējams, uzklikšķinot uz pluszīmes rindīņas sākumā.

Novietojot peli uz jebkuras jūs interesējošā vārda ailītes markējuma tabulā (ja vārdam jau ir pievienots morfoloģiskais raksturojums) (sk. 4. att.) vai izvēlnes daļā, kur doti iespējamie analīzes varianti (sk. 5. att.), papildu logā automātiski parādās izvēsta informācija par šo vārdu (galvenokārt morfoloģiskais raksturojums).

Markētājs - 0.9.59 [C:\Users\Iļze Auziņa\Desktop\annotetajs 22.05.09\annotator-r592\demo1.txt]

Korpuss Rīki Palīgs

Čānks

Markējuma tabula Sintakse

Dabiski, ka ūdens nevar pārvērsties par zivi vai taureni.

Ūdens vispār nevar pārvērsties.

Tīrs ūdens mūžīgi ir tīrs ūdens.

Tātad Parmenīdam bija taisnība, ka "nekas nepārvēršas".

Tajā pašā laikā Empedokls bija vienprātis ar Heraklītu, ka mums jāpajaujas.

Mums jātic tam, ko redzam, un mēs redzam tieši to, ka.

dabā viss pastāvīgi mainās.

Empedokls secināja, ka jāatmet pieņēmums par vienas pirmvielas pastāvēšanu.

Ne ūdens, ne gaiss paši par sevi nevar pārvērsties par rožu krūmu vai taureni.

Tātad nav iespējams, ka dabā ir tikai viena viela, kas ir visa pamatā.

Empedokls apgalvoja, ka dabā pastāv pavisam četri pamatelementi jeb "cēloņi".

Tie, viņaprāt, bija zeme, gaiss, uguns un ūdens.

Visas pārvērtības dabā nāca no tā, ka šie četri elementi savienojās un atkal sadalījās.

Jo viss sastāv no zemes, gaisa, uguns un ūdens, tikai dažādās attiecībās.

Kad zieds novīst vai kāds

✓ Rādīt visus vārdus

Npk.	Vārds	Markējums	Pamatforma	Loma
1	Dabiski	r_	dabiski	
2	,	zc	,	
3	ka	css	ka	
4	ūdens	ncmsn2	ūdens	
5	nevar			
6	pārvērsties			
7	par			
8	zivi			
9	vai			
10	taureni			
11	.			

Vārds	Markējums	G...	Ti...	Apraksts
+ nevar	<[v,o,n,i,p,t,3,3,0,a,y], 'nevar'>	474	173	Darbības
+ nevar	Es zinu labāk!			
+ nevar	Reziduālis			

Vārdšķira = Lietvārds
Galotnes nr = 39
Vārds = ūdens
Deklinācija = 2
Mija = 0
Locījums = Nominatīvs
Pamatforma = ūdens
Vārdgrupas nr = 4
Leksēmas nr = 8749
Skaitlis = Vienskaitlis
Lietvārda tips = Sugas vārds
Dzimte = Vīriešu
Avots = Valērija leksikons
Kristīnes markējums = ncmsn2 adne, 16, 21

4. att. Papildinformācija par vārdformu *ūdens*

Empedokls secināja, ka jāatmet pieņēmums par vienas pirmviela	6	pārvērsties	vmyn0l1000n	pārvērsties		
Ne ūdens, ne gaiss paši par sevi nevar pārvērsties par rožu krā	7	par				
Tātad nav iespējams, ka dabā ir tikai viena viela, kas ir visa par	8	zivi				
Empedokls apgalvoja, ka dabā pastāv pavisam četri pamatelem	9	vai				
Tie, viņaprāt, bija zeme, gaiss, uguns un ūdens.	10	taureni				
Visas pārvērtības dabā nāca no tā, ka šie četri elementi savien	11	.				
Jo viss sastāv no zemes, gaisa, uguns un ūdens, tikai dažādās						
Kad zieds novīst vai kāds						
<input type="checkbox"/> Rādīt visus vārdus						
Vārds	Markējums	G...	Ti...	Apraksts		
+ par	<[s,p,p,d,[n]],'par'>	1196	1	1196, 25, Prievārds, par, 28253, Daudzskaitlis, Datīvs, 0, Pirms, Nē, par, sppdn		
+ par	<[s,p,s,a,[n]],'par'>	1196	1	1196, 25, Prievārds, par, 28254, Vienskaitlis, Akuzatīvs, 0, Pirms, Nē, par, sspan		
+ par	Es zinu labāk!					
+ par	Reziduālis					

Galotnes nr = 1196
Vārdgrupas nr = 25
Vārdšķira = Prievārds
Vārds = par
Leksēmas nr = 28253
Skaitlis = Daudzskaitlis
Rekcija = Datīvs
Mija = 0
Novietojums = Pirms
Vietas apstākļa nozīme = Nē
Pamatforma = par
Kristīnes markējums = sppdn

5. att. Papildinformācija par vārdu *par*

Pareizā markējuma izvēle

Latviešu valodas markētājs izmanto analizatora vārdnīcā (leksikonā) atrodamo informāciju par vārdiem. Ja kāds no teikumā esošajiem vārdiem neatrodas analizatora vārdnīcā, tātad nav iepriekš zināma tā vārdšķira u. c. informācija par vārdu, analizators mēģina to noteikt pēc vārdformas morfoloģiskajām pazīmēm (galotnes, piedēkļa, piedēkļa).

Tāpēc daudzos gadījumos programma piedāvā vairākus markējuma variantus, jo, kā zināms, latviešu valodā diezgan bieži ir sastopamas homoformas, t. i., divas vai vairākas dažādas nozīmes vārdformas, kuras raksta vienādi un kuras nav iespējams atšķirt pēc ārējām pazīmēm (bez leksiskās nozīmes analīzes), piemēram, *sviesta* (lietvārds *sviests* vienskaitļa ģenitīvā) un *sviesta* (divdabis *sviests* sieviešu dzimtes vienskaitļa nominatīvā), *kur* (apstākļa vārds) un *kur* (darbības vārda *kurt* īstenības izteiksmes vienkāršās tagadnes vienskaitļa 2. personas un vienskaitļa un daudzskaitļa 3. personas forma), *mēs* (daudzskaitļa 1. personas vietniekvārds) un *mēs* (darbības vārda *mēt* īstenības izteiksmes vienkāršās nākotnes 3. personas forma) (sk. 6. att.). Markētājs piedāvā risinājumu homoformu analīzē: lietotājam, uzklikšķinot uz pluszīmes, ir iespēja izvēlēties atbilstošo morfoloģisko raksturojumu.

Mēs to varam salīdzināt ar gleznošanu.		Npk.	Vārds	Marķējums	Pamatforma	Loma	Paskaidro
		1	Mēs			teikuma priek...	izteicējs
		2	to			tiešais papild...	izteicējs
		3	izteicējs	vmnpt01pan		izteicējs	.
		4	varam				izteicējs
		5	salīdzināt				izteicējs
		6	pievārdeklis	n_fsa_		netiešais papi...	izteicējs
		7	ar				pievārdeklis
		8	gleznošanu				pievārdeklis
		9	.				
		<input type="checkbox"/> Rādīt visus vārdus					

	Vārds	Marķējums	Gal...	Ticamība	Apraksts
+	Mēs	<[v,m,n,i,f,i,1,3,0,a,n],mēt>	717		1 717, Darbības vārds, 1, mēs, 6, Nākotne, mēt, 14, Nepiemīt, 24074, Nē, Īstenības, Nepāre...
+	Mēs	<[v,n,i,f,i,2,3,0,a,n],mēt>	243		1 243, Darbības vārds, 2, mēs, 0, Nākotne, mēt, 15, Nepiemīt, 46859, Nē, Īstenības, 3, Valē...
+	Mēs	<[p,p,1,0,p,n,_[0]],mēs>	2011		1 2011, Vietniekvārds, mēs, Nepiemīt, 0, Nominatīvs, Personu, mēs, jā, 24, Daudzskaitis, 27...
+	Mēs	Es zinu labāk!			
+	Mēs	Reziduālis			

6. att. Vārdformas *mēs* morfoloģiskās analīzes varianti

Opcija *Es zinu labāk!*

Ja neviens no vārdformas raksturojumiem nav pareizs, programmas lietotājs var piedāvāt savu vārda morfoloģisko raksturojumu, uzklikšķinot uz *Es zinu labāk!*, un vienlaicīgi to pievienot jau eksistējošam leksikonam. Arī tad, ja informācijas par kādu vārdu vārdnīcā nav un marķētājs nepiedāvā nevienu analīzes variantu, lietotājs to var pievienot pats, uzklikšķinot uz *Es zinu labāk!*

- Parādīsies logs (sk. 7. att.), kura ailītē *Pamatforma* jāieraksta konkrētā (analizējamā) vārda pamatforma.
- Pēc pamatformas ierakstīšanas, ailē *Vārdšķiras* automātiski tiks piedāvātas iespējamās vārdšķiras, arī lietvārdu deklinācijas, verbu konjugācijas u. tml.
- Tālākā izvēle atkarīga no vārda piederības konkrētai vārdšķirai, t. i., ja analizējamais vārds ir darbības vārds, būs iespējams rakstot 1) verba transitivitāti (pārejošs/nepārejošs), 2) noteikt darbības vārda tipu (*Patstāvīgs darbības vārds/ palīgdarbības vārds/ Modāls/ Fāzes/ Izpaušmes veida/ Palīgverbs 'būt'/ Palīgverbi 'tikt' un 'tapt'*). Ja analizējamais vārds ir lietvārds, var raksturot tā tipu, t. i., vai tas ir īpašvārds vai sugas vārds.
- Ja analizētajam vārdam nepieciešamās pazīmes ir pievienotas, tiek piedāvāta iespēja [Pievienot vārdu] vai [Atcelt pievienošanu]. Izvēloties iespēju pievienot vārdu, tas ar attiecīgi izveidoto morfoloģisko marķējuma variantu parādīsies saskarnes lejasdaļā, kur to būs iespējams apstiprināt, nospiežot pluszīmi.

7. att. Vārda pievienošana leksikonam

Ja lietotājs nav apmierināts ar piedāvāto morfoloģisko vārdformas raksturojumu, to var labot manuāli, uzklikšķinot uz atbilstošās ailītes saskarnes labajā pusē slejā *Marķējums*.

Lai redzētu visas jau samarķētās teksta vienības, saskarnes logā (zem vārdformu un pieturzīmju saraksta) jāatzīmē *Rādīt visus vārdus*.

Marķētājs piedāvā iespēju pārmarķēt jau apstrādātu tekstu vai arī turpināt marķēt citā teksta vietā, piemēram, izlaižot vairākus teikumus. Lai to izdarītu, jāiezīmē rinda (teksta fragments), no kuras teksts tiks pārmarķēts, un ar peles labo taustiņu jāuzklikšķina uz tās. Tiks piedāvāta izvēlne, kurā jāizraugās iespēja [Turpināt marķēt no šīs rindas].

Morfoloģiskās pazīmes

Morfoloģiskās analīzes laukā katrai vārdformai tiek noteikta morfoloģisko pazīmju kopa. Pazīmju skaits kopā katrai vārdšķīrai ir atšķirīgs, piemēram, darbības vārds tiek raksturots, izmantojot 11 pazīmes (sk. 8. att.), lietvārds — izmantojot 6 pazīmes, savukārt izsaukmes vārda un partikulas raksturošanai izmanto tikai 2 pazīmes.

Pirmā pazīme (kodifikators) morfoloģisko pazīmju kopā norāda vārda piederību konkrētai vārdšķīrai, piemēram, darbības vārda kodifikators ir *v*, lietvārda — *n*, īpašības vārda — *a*, vietniekvārda — *p*, apstākļa vārda — *r*.

Vairākas pazīmes dažādu vārdšķīru vārdiem ir kopīgas, piemēram, dzimte, skaitlis, locījums piemīt lietvārdam un īpašības vārdam (3., 4., 5. pazīme), vietniekvārdam (4., 5., 6. pazīme), divdabim (6., 7., 8. pazīme).

Atsevišķos gadījumos morfoloģisko pazīmju kopā iekļautas arī semantiskas pazīmes, piemēram, kur iespējams, apstākļa vārdiem norādīts, kādas nozīmes apstākļa vārdi tie ir — laika, vietas, veida, mēra, cēloņa. Saikļi raksturoti gan pēc sintaktiskās funkcijas, gan pēc uzbūves, piemēram, saiklis *jo* — vienkāršs pakārtojuma saiklis un atkārtots sakārtojuma saiklis. Norādāmajiem vietniekvārdiem atzīmēts, kādas vārdšķīras tie aizstāj.

□[v,m,n,i,p,t,3,3,p,a,n], □ liecināt'>

Pazīme	vērtība	atsifrējums
kodifikators	v	darbības vārds
tips	m	patstāvīgs (<i>main</i>)
atgriezeniskums	n	nē
izteiksme	i	īstenības izteiksme
laiks	p	tagadne (<i>present</i>)
pabeigtība	i	nepabeigts (<i>imperfect</i>)
konjugācija	3	3. konjugācija
persona	3	3. persona
skaitlis	p	daudzskaitlis (<i>plural</i>)
kārta	a	darāmā kārta
noliegums	n	nav noliegtais darbības vārds

8. att. Vārdformas *liecina* morfoloģiskais raksturojums

Sintaktiskā analīze

Ja iespējams, sintaktiskā informācija marķētajam teksta fragmentam (konkrēti — vārdformām), tiek pievienota automātisku. Pamatojoties uz vārdu morfoloģiskajām formām, tiek analizētas teikuma sintaktiskās attiecības (piem., saskaņojums, atkarības), noteikts, piemēram, ka lietvārds akuzatīvā var būt atkarīgs no verba finītās formas, kā arī tiek noteikta vārdformas sintaktiskā loma teikumā. Analizētājs analizē teikumu gramatiski (morfoloģiski un sintaktiski) analīzi — nosaka katra vārda morfoloģisko formu un, balstoties uz tām, parāda iespējamo teikuma sintaktisko struktūru.

Lai vārdformu analīzes tabulai pievienotu sintaktiskās lomas, atkal jāiezīmē konkrēts teikums/teksta fragments un jāuzklikšķina ar peles labo taustiņu. Parādoties izvēlnei, jāizraugās [Izvēlēties citu analīzes variantu]. Tad jaunā logā tiks piedāvāts teikuma sintaktiskais marķējums, t. s. kastītes (sk. 9. att.).

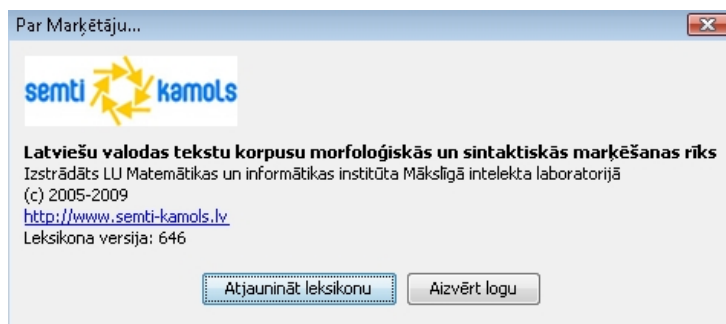


9. att. Sintaktiskās analīzes piemērs

Ja tiek piedāvāti vairāki varianti, lietotājam jāizraugās visatbilstošākais, uzklikšķinot uz saskarnes labajā malā redzamā cipara. Tad automātiski jaunā informācija par vārdu sintaktiskajām lomām tiks ievietota iepriekš aprakstītās tabulas (sk 1. att., 7. att.) 5. un 6. ailē. Ja neviens no piedāvātajiem variantiem nav precīzs / pareizs, jānospiež saskarnes lejā redzamā poga *Neviens no variantiem īsti neatbilst*.

Leksikona atjaunināšana

Marķējot tekstus, ik pa laikam leksikonam tiek pievienots jauns vārds (sk. 7. att.). Tā kā marķētāju vienlaicīgi var izmantot vairāki lietotāji, kas savukārt neatkarīgi cits no cita leksikonam var pievieno vārdus, leksikons ir jāatjaunina. To var izdarīt, ejot uz izvēlni [Palīgs] \ [Par...] un izvēlnes logā, uzklikšķinot uz pogas [Atjaunināt leksikonu] (sk. 10. att.). Pašreizējā versijā lejupielādētais atjauninājums ir pusautomātiski jāsinhronizē ar esošo leksikonu, izmantojot [Rīki] \ [Sinhronizēt leksikonus].



10. att. Leksikona atjaunināšana