

Eeltöötlusmooduli dokumentatsioon

Eesmärk: süntaksianalüüsi-eelse eeltöötlusmoodul parandab korpustes lausestusvigu ja sagedasemaid süntaktilise analüüsi vigu põhjustavad vead. Lausestusvigadest parandatatakse selliseid juhtumeid, kus lausemärgendid on üleliigsed. Süntaktilise analüüsi vigu põhjustavate vigade korral ühendati lause sees paiknevad mitteleksikaalsed elemendid täiendava märgendiga (<+>). Spordiuudiste ja saatekava lõikude puhul ignoreeritakse mitte-lauselisi sisu.

1. Arvud ja numbrid:

- $20\ 000 \Rightarrow 20<+>000$
- $669\ 81\ 54 \Rightarrow 669<+>81<+>54$
- $1\ ,\ 2 \Rightarrow 1,2$
- $25\ -\ 30\ \% \Rightarrow 25-30\% ,\ 20\ \%ga \Rightarrow 20\%ga, 0,20\ -protsendilise \Rightarrow 0,20-protsendilise$
- lühend (nt AS) või jutumärk, kaldkriips, § või arv, millele järgneb käändelõpp:
 - $1,5\ -ni \Rightarrow 1,5-ni$
 - $\S\ -st \Rightarrow \S-st$
- avaldised:
 - $24+9 = 33 \Rightarrow 24+9<+>=<+>33$
 - $R\ 2 = 0\ ,\ 08 \Rightarrow R<+>2<+>=<+>0,08$
 - $r = 0\ ,\ 46 \Rightarrow r<+>=<+>0,46$
- arv, millele järgneb mõõtühiku lühend:
 - $2\ t\ 15\ min \Rightarrow 2<+>t\ 15<+>min$
 - $294\ ha-lt \Rightarrow 294<+>ha-lt$
 - $1\ ha \Rightarrow 1<+>ha$
 - $1,0\ mM \Rightarrow 1,0<+>mM$
 - $740\ kHz-ni \Rightarrow 740<+>kHz-ni$
 - $7\ C^{\circ} \Rightarrow 7<+>C^{\circ}$
 - $60\ km\ /\ h \Rightarrow 60<+>km/h$
 - $2,3\ h\ /\ m \Rightarrow 2,3<+>h/m$
 - $443,49\ kr/MWh \Rightarrow 443,49<+>kr/MWh$
 - $22\ 000\ kr/m^2 \Rightarrow 22<+>000<+>kr/m^2$
- aastarvud ja kuupäevad:
 - $21.\ </s>\ <s>\ 12.\ 2001 \Rightarrow 21.12.2001$
 - $1884.\ a. \Rightarrow 1884.a.$
 - $1884\ a. \Rightarrow 1884a.$
 - $15.\ 04.\ 2005 \Rightarrow 15.04.2005$
 - $1945\ .\ aasta \Rightarrow 1945.\ aasta$

2. Vahemikud:

- $40\ 000-45\ 000 \Rightarrow 40<+>000<+>-<+>45<+>000$
- $0\ ,\ 3\ \dots\ </s>\ <s>\ 1\% \Rightarrow 0,3<+>...<+>1\%$

- 1998. - 2000 \Rightarrow 1998.<+>-<+>2000
- 2,0 ... </s> <s> 3,5 \Rightarrow 2,0<+>...<+>3,5
- 40 – 300 C° \Rightarrow 40<+>-<+>300<+>C°
- 3 ... 8 mÜs \Rightarrow 3<+>...<+>8<+>mÜs
- 4 . - 5 . \Rightarrow 4.<+>-<+>5.

3. Nimed:

- J. Fr . </s> <s> Blumenbach \Rightarrow J.Fr.<+>Blumenbach
- A . </s> <s> J. Sjögren \Rightarrow A.J.<+>Sjögren
- J. R. R. Tolkieni \Rightarrow J.R.R.<+>Tolkieni
- St. Louis \Rightarrow St.<+>Louis
- Simon & Schusteri \Rightarrow Simon<+>&<+>Schusteri
- Laulasmaa Spa & Konverentsihotell \Rightarrow Laulasmaa Spa<+>&<+> Konverentsihotell

4. Lausestusvead:

- (ala)jaotused, loetelud, järgarvud:
 - (2 . </s> <s> 8) \Rightarrow (2 . 8)
 - (85 . </s> <s> Antonov) \Rightarrow (85 . Antonov)
 - 1 . </s> <s> 4. 3. </s> <s> Tahtlus /.../ \Rightarrow 1.4.3. Tahtlus /.../
 - 5 . </s> <s> 2. 3. </s> <s> Ujumisõpetajad \Rightarrow 5.2.3.Ujumisõpetajad
- veebiaadressid:
 - www. </s> <s> RODEsign.ee/ \Rightarrow www.RODesign.ee/
- lauselõpud, kus punkt on stringi külge jäänud:
 - ajal. </s> Peame /.../ \Rightarrow ajal . </s> Peame /.../
 - jne. </s> <s> Aga udmurdid /.../ \Rightarrow jne . </s> <s> Aga udmurdid /.../
 - aastasse 2000. </s> <s> Ja teine /.../ \Rightarrow aastasse 2000 . </s> <s> Ja teine /.../

5. Spordiuudised:

- tulemused, ajad:
 - 7 : 8 \Rightarrow 7<+>:<+>8
 - 2 . 06 , 08 \Rightarrow 2<+>.<+>06<+>,<+>08

6. <ignore> ja </ignore> märgendid:

- sulgudes olevad arvud ja/või lühendid, suurtähega algavad sõned:
 - (WTA 210.) \Rightarrow <ignore> (WTA 210.) </ignore>
 - Kreekaga (57.) \Rightarrow Kreekaga <ignore> (57.) </ignore>
 - Tuneesia (5) \Rightarrow Tuneesia <ignore> (5) </ignore>
- pikad loetelud:

- <p> <s> Rahvamajandus : Abiševa , Maria ; Ahlamtsenkova , Viktoria ; /.../ Vinkel , Natalja ; Zahharov , Sergei . </s> </p> ⇒ <p> <ignore> <s> Rahvamajandus : Abiševa , Maria ; Ahlamtsenkova , Viktoria ; Bakulina , Maria ; /.../ Vinkel , Natalja ; Zahharov , Sergei . </s> </ignore> </p>
- <p> <s> Väravad : Vahtramäe (5) , Ustritski (7) , M. Rooba (34) - Voronin (30) , Leitan (42) </s> </p> ⇒ <p> <ignore> <s> Väravad : Vahtramäe (5) , Ustritski (7) , M.<+>Rooba (34) - Voronin (30) , Leitan (42) </s> </ignore> </p>
- <p> <s> Tabeliseis : Austria 6 punkti , Poola 4 , Leedu ja Holland 3 , Eesti 2 , Horvaatia 0 . </s> </p> ⇒ <p> <ignore> <s> Tabeliseis : Austria 6 punkti , Poola 4 , Leedu ja Holland 3 , Eesti 2 , Horvaatia 0 . </s> </ignore> </p>
- spordiuudised:
 - lõigud või laused, mis algavad spordiala ja/või distantsiga ning millele järgneb tulemuste loetelu:
 - * <p> <s> 5 km (v) : 1. Katerina Neumannova (pildil , Tshehhimaa) 12.56 , 1 , 2. </s> <s> Tshepalova +9,0 , 3. </s> /.../ <s> Jelena Buruhina (Venemaa) +39,2 , 28. </s> <s> Skari +53,8 , ... </s> <s> 65. Katrin *migun +1.33 , 0. </s> </p> ⇒ <p> <ignore> <s> 5 km (v) : 1. Katerina Neumannova (pildil , Tshehhimaa) 12.56 , 1 , 2. </s> <s> Tshepalova +9,0 , 3. </s> /.../ <s> Jelena Buruhina (Venemaa) +39,2 , 28. </s> <s> Skari +53,8 , ... </s> <s> 65. Katrin *migun +1.33 , 0. </s> </ignore> </p>
 - * <p> <s> Sprint : 1. Ronny Ackermann (Saksamaa) 18.58 , 2 , 2. </s> <s> Hannu Manninen (Soome) +9,3 , 3. </s> <s> Kristian Hammer (Norra) +10,4 , 4. </s> /.../ <s> Christoph Bieler (Austria) +1.09 , 1 , 25. </s> <s> Jens Salumäe (Eesti) +1.31 , 0. </s> </p> ⇒ <p> <ignore> <s> Sprint : 1. Ronny Ackermann (Saksamaa) 18.58 , 2 , 2. </s> <s> Hannu Manninen (Soome) +9,3 , 3. </s> <s> Kristian Hammer (Norra) +10,4 , 4. </s> /.../ <s> Christoph Bieler (Austria) +1.09 , 1 , 25. </s> <s> Jens Salumäe (Eesti) +1.31 , 0. </s> </ignore> </p>
 - * <p> <s> Tulemused . </s> <s> Naiste 1000 m : 1 . </s> <s> Chris Witty (USA) 1 . 14 , 96 , 2 . </s> /.../ <s> Elena Belci (Itaalia) 7 . 16 , . </s> </p> ⇒ <p> <s> Tulemused . </s> <ignore> <s> Naiste 1000<+>m : 1 . </s> <s> Chris Witty (USA) 1<+>. <+>14<+>,<+>96 , 2 . </s> /.../ <s> Elena Belci (Itaalia) 7<+>. <+>16<+>,<+>43 . </s> </ignore> </p>
 - * <p> <s> Mehed : 200 m 1 . </s> <s> Maurice Greene (USA) 19 , 92 , 2 . </s> <s> Frank Fredericks (Namibia) 19 , 93 , 3 . </s> /.../ <s> Younes Moudrik (Maroko) 8 . 20 , 3 . </s> <s> Kareem Streete-Thompson (Kaimanisaared) 8 . 15 . </s> </p> ⇒ <p> <ignore> <s> Mehed : 200<+>m 1 . </s> <s> Maurice Greene (USA) 19 , 92 , 2 . </s> <s> Frank Fredericks (Namibia) 19 , 93 , 3 . </s> /.../ <s> Younes Moudrik (Maroko) 8<+>. <+>20<+>,<+>3 . </s> <s> Kareem Streete-Thompson (Kaimanisaared) 8 . 15 . </s> </ignore> </p>
 - * <p> <s> <hi rend="rasvane"> 10 000 m : </hi> 1 . </s> <s> Sally Barsosio (Keenia) 31 . 32 , 29 , 2 . </s> /.../ <s> Chemi Takahashi (Jaapan) 32 . 23 , >61 . </s> </p> ⇒ <p> <ignore> <s> <hi rend="rasvane"> 10<+>000<+>m : </hi> 1 . </s> <s> Sally Barsosio (Keenia) 31<+>. <+>32<+>,<+>29 , 2 . </s> /.../ <s> Chemi

- Takahashi (Jaapan) 32<+>.<+>23<+>,<+>61 . </s> </ignore>
</p></div4>
- * <p> <s> <hi rend="rasvane"> 400 m : </hi> 1 . </s> <s> Iwan Thomas (Suurbritannia) 44 , 90 , 2 . </s> /.../ <s> Chris Jones (USA) 45 , 71 , 6 . </s> <s> Sunday Bada (Nigeeria) 45 , 86 . </s> </p> ⇒ <p> <ignore> <s> <hi rend="rasvane"> 400<+>m : </hi> 1 . </s> <s> Iwan Thomas (Suurbritannia) 44 , 90 , 2 . </s> /.../ <s> Chris Jones (USA) 45 , 71 , 6 . </s> <s> Sunday Bada (Nigeeria) 45 , 86 . </s> </ignore> </p>
 - * <p> <s> <hi rend="rasvane"> Maraton </hi> 2.25.17 Rosa Mota , Portugal 1987 </s> </p> ⇒ <p> <ignore> <s> <hi rend="rasvane"> Maraton </hi> 2.25.17 Rosa Mota , Portugal 1987 </s> </ignore> </p>
- meeskonna või sportlase nimi, millele järgnevad tulemused (numbrijada):
- * <p> <s> Joe Sakic Col 32 16+27= 43 10 </s> </p> ⇒ <p> <ignore> <s> Joe Sakic Col 32<+>16+27=<+>43<+>10 </s> </ignore> </p>
 - * <p> <s> NY Islanders 20 6 3 0 11 15 </s> </p> ⇒ <p> <ignore> <s> NY Islanders 20<+>6<+>3<+>0 11<+>15 </s> </ignore> </p>
 - * <p rend="rasvane"> <s> Miami-Orlando 2 : 2 </s> </p> ⇒ <p rend="rasvane"> <ignore> <s> Miami-Orlando 2<+>:<+>2 </s> </ignore> </p>
 - * <p> <s> NEW JERSEY 14 29 . </s> <s> 326 18 1/2 </s> </p> ⇒ <p> <ignore> <s> NEW JERSEY 14<+>29 . </s> <s> 326<+>18<+>1/2 </s> </ignore> </p>
 - * <p> <s> Zimbru Chisinau (Moldova) – Dinamo Tbilisi (Gruusia) 2 : 0 (kokku 3 : 2) . </s> /.../ <s> Dariusz Gesior 14 , Artur Wichniarek 52 , 60 , Radoslaw Michalski 75 – Svetoslav Todorov 30 . </s> </p> ⇒ <p> <ignore> <s> Zimbru Chisinau (Moldova) – Dinamo Tbilisi (Gruusia) 2<+>:<+>0 (kokku 3<+>:<+>2) . </s> /.../ <s> Dariusz Gesior 14 , Artur Wichniarek 52 , 60 , Radoslaw Michalski 75 – Svetoslav Todorov 30 . </s> </ignore> </p>
 - * <p> <s> Teised kohtumised : Washington - Colorado 2 : 1 , Ottawa - Carolina 5 : 1 , Tampa Bay - NY Rangers 2 : 2 (la) , NY Islanders - San Jose 1 : 4 , St . </s> <s> Louis - Atlanta 4 : 1 , Nashville - Edmonton 1 : >2 . </s> </p> ⇒ <p> <ignore> <s> Teised kohtumised : Washington - Colorado 2<+>:<+>1 , Ottawa - Carolina 5<+>:<+>1 , Tampa Bay - NY Rangers 2<+>:<+>2 (la) , NY Islanders - San Jose 1<+>:<+>4 , St . </s> <s> Louis - Atlanta 4<+>:<+>1 , Nashville - Edmonton 1<+>:<+>2 . </s> </ignore> </p>
 - * <p> <s> 2. Svetlana Tšernoussova Venemaa +1.12 , 8 (1) </s> </p> ⇒ <p> <ignore> <s> 2. Svetlana Tšernoussova Venemaa +1.12 , 8 (1) </s> </ignore> </p>
 - * <p> <s> 1. Valencia 17 10 5 2 29 : 10 35 </s> </p> ⇒ <p> <ignore> <s> 1. Valencia 17<+>10<+>5<+>229<+>:<+>10<+>35 </s> </ignore> </p>
 - * <p> <s> K Clijsters (BEL) (15) - A Jidkova (RUS) 6 - 3 , 7 - 6 </s> </p> ⇒ <p> <ignore> <s> K Clijsters (BEL) (15) - A Jidkova (RUS) 6<+>-<+>3 , 7<+>-<+>6 </s> </ignore> </p>
- viited teadustekstides (vt PATT_BRACS, PATT_12):

- (vt joonis 4.6) ⇒ <ignore> (vt joonis 4.6) </ignore>
- (2001 : 114) ⇒ <ignore> (2001 : 114) </ignore>
- (Silver , 1992 ; Alessi& Trolip , 2001 : 115) ⇒ <ignore> (Silver , 1992 ; Alessi& Trolip , 2001 : 115) </ignore>
- (1984) ⇒ <ignore> (1984) </ignore>
- (Miljan , R. et al. Choices for Dairy ... 2003) ⇒ <ignore> (Miljan , R. et al. Choices for Dairy ... 2003) </ignore>
- (tabel 3.2) ⇒ <ignore> (tabel 3.2) </ignore>
- sulgude sees paiknevad lühendid, sõnede ja arvude kombinatsioonid (vt PATT_BRACS, PATT_12):
 - (A - 23,87 tuhat krooni , B - 23,72 tuhat krooni ja D - 23,35 tuhat krooni) ⇒ <ignore> (A - 23,87 tuhat krooni , B - 23,72 tuhat krooni ja D - 23,35 tuhat krooni) </ignore>
 - (195 miljonit USDd) ⇒ ignore> (195 miljonit USDd) </ignore>
 - (à 30 min) ⇒ <ignore> (à 30<+>min) </ignore>
 - (“ Hannah ja ta õed ” , 1986) ⇒ <ignore> (“ Hannah ja ta õed ” , 1986) </ignore>
 - (Hispaania , 3) ⇒ <ignore> (Hispaania , 3) </ignore>
 - (rasvasus - 12% , valgusisaldus - 16%) ⇒ <ignore> (rasvasus - 12% , valgusisaldus - 16%) </ignore>
 - (2Kr 6,10) ⇒ <ignore> (2Kr 6,10) </ignore>
 - (15 juhtumit , kusjuures 2 said surma) ⇒ <ignore> (15 juhtumit , kusjuures 2 said surma) </ignore>
 - (50 : 47,81 : 81 ; Jamchy 20 , Henefeld 19 - Gadou 25 , Rigauveau 24) ⇒ <ignore> (50 : 47,81 : 81 ; Jamchy 20 , Henefeld 1 9<+>-<+>Gadou 25 , Rigauveau 24) </ignore>
 - (RKO) ⇒ <ignore> (RKO) </ignore>
 - (viimasel päeval võitis Goran Ivanisevic (Horvaatia) Thomas Musteri 6 : 7 , 7 : 5 , 6 : 7 , 6 : 2 , 7 : 5) , Zimbabwe-Suurbritannia 4 : 1 . </s> </p> ⇒ <ignore> (viimasel päeval võitis Goran Ivanisevic (Horvaatia) Thomas Musteri 6<+>:<+>7,7<+>:<+>5,6<+>:<+>7,6<+>:<+>2,7<+>:<+>5) </ignore> , Zimbabwe-Suurbritannia 4<+>:<+>1 . </s> </p>
 - (100 ja 4x100 m) ⇒ <ignore> (100 ja 4x100 m) </ignore>
 - <s> 26-aastase Aleksander Tammerti seeria (65.36 - 59.44 - 66.95 - 62.06 - 62.50 - 66.48) andis /.../ ⇒ <s> 26-aastase Aleksander Tammerti seeria <ignore> (65.36-59.44-66.95-62.06-62.50-66.48) </ignore> andis /.../
- telekava:
 - <p> <s> <hi rend="rasvane"> 06 . 00 </hi> Saatekavateade <hi rend="rasvane"> 06 . 05 </hi> /.../ Õigel ajal <hi rend="rasvane"> 23 . 05 </hi> Venemaa aidi vastu <hi rend="rasvane"> 23 . 55 </hi> Ilmateade <hi rend="rasvane"> 00 . 00 </hi> Kunstipoodium <hi rend="rasvane"> 00 . 35 </hi> Telepood </s> </p> ⇒ <p> <ignore> <s> <hi rend="rasvane"> 06 . 00 </hi> Saatekavateade <hi rend="rasvane"> 06 . 05 </hi> /.../ Õigel ajal <hi rend="rasvane"> 23 . 05 </hi> Venemaa aidi vastu <hi rend="rasvane"> 23 . 55 </hi> Ilmateade <hi rend="rasvane"> 00 . 00 </hi> Kunstipoodium <hi rend="rasvane"> 00 . 35 </hi> Telepood </s> </ignore> </p>

– <p> <s> <hi rend="rasvane"> 07 . 00 </hi> Tere hommikust ! </s> </p>
 ⇒ <p> <ignore> <s> <hi rend="rasvane"> 07 . 00 </hi> Tere hommikust !
 </s> </ignore> </p>

7. Igale lausele lisati selle järjekorranumber tekstis (va <ignore> ja </ignore> märgendite vahel olevatele lausetele):

- <p> <s> <id="15">
- </s> <s> <id="1066">