

Ðàçìå÷åííúé êîðïóñ äèàëíáíâ êàê Õåñóðñ ìíäåéèðíâàíéý äèàëíáà: ïðãàíèçàöèý è ðàçìåòéà Ý

Àâòíð ìéññá Áéññáðíá Íññéááéáíâ

18.07.2008 á.

Íññéááéáíâ 20.07.2008 á.

Òåêñòú íå÷àòííáí âçääàíéý â ôíðìàòå PDF

Ýñòíññéé êîðïóñ äèàëíáíâ ñíçääáí ëeíññéñòáíè è êíññüþòåðíúè òåññéíáàíè Òàðòóññéíáí óíèååðñèòåòà ñ öåéüþ èñññéåäåàòü ðå-
÷åéíáíâ÷åñéíé êíññéíééè.

1. Niññòáâ êîðïóñà

Íññéááéáíâ ÷àñòü êîðïóñà ñíññòíèò èç åñòåñòååííúô óñòíúô äèàëíáíâ (758 ñíðàâí÷íúô òåéååðííúô äèàëíáíâ è 106 íáññòååñòååííúô ä-

Óñòíúâ äèàëíáíè ñíññðåíýþòñý â öèôðíâííí èéè ìøèôðíâàííí õíðìàòå .wav, à â òåêñòíâíí áèäå ïðåäñòååéåíû â ðàíññéðèíøè Áæåôô
Òðàíññéðèíøè ðèñéðóà ÿåéåíéý, íçâíéþþùèå ïðíññéååèòü äèíàíè÷åñéíâ ïññòðíâíèå äèàëíáà èç ðåéèé è èíòíàðèííúô ååéíéë:

Íàëäý, íí áàæíàý ÷àñòü êîðïóñà —; 22 íèññìáííúô äèàëíáà íáúáíí â 2500 ñéíâ —; ïðåäñòååéýåò ñíáíé êíññüþòåðíúâ ñí
åñòåñòååíííí ÿçûêå, íí â äåéñòåèòåéüññòé ðíëü íðíñðàííü èñññéíýé ÷åéíååé. Äèàëíáè ïðåäñòååéýþò ñíáíé eíäé áåñåå ííëüçíâàòå

2. Óðíâíè àíññòåöè

2.1. Íðôíëíáè

Хàñòü óñòíúô äèàëíáíâ êîðïóñà (20000 ñéíâ) ðàçìå÷åíà ííðôíëíáè÷åñéíé. Àâòíðàòè÷åñéàý ðàçìåòéà ïñóùåñòåéåíâ ñ ííññùüþ ííðôíëíáè
íññóæååííéå, êíòíðííó ñíññòåñòååíâéí èñññéååíâàíéå íðíñðàííéí ñíýðèý ííðôíëíáè÷åñéíé ííññéíé è óñòííé ýñðíññéíé ðå÷é â ñðåáåíáíèè ñ

2.2. Ðå÷åååå åàêòú

Íe iāia èç ðàññiñòðåíûô ñoñàì ðàçìåòèè íá ìòâå÷àëà öäëëi ïiääëëðiâàíëy è èññëåäiâàíëy áñòåñòâåíûô äëàëïâàí ñ òi÷êè çðåíëy òëiñëiæé [Gerassimenko et al. 2004].

Òèiñëiæy ñiñòièò èç 126 äèàëiââûô àêòiâ è ïñiââàåòñy íà ðàññiñòðåíèè äèàëiâà èàê ñiöèàëüíâi âçàëiñäéñòâèy. Ðå÷ââûâàèiñòiñûâ àêòôû (íáiñòiââû, iñáóëèòåëüñûâ ðâiñëèéè íåðâòiñé ñâýçè). Ñ äðóâié ñòiðñûâ, àêòôû iñáðâçäåéýþòñy íá ìòíñyùeâñy ê iðåâiñûâàòåëüñé iðiñâðâiñûâ.

Íðèiâð 1 (âñòåñòâåíûé äèàëiâà íà ìâðåââââå íà ðóññêèé).

((çâiñîê)) | RIE: ÇÂIñîÊ |

V: `ñiñòaâiñòiâëy | RIJ: ÍÒÅÅÒ ÍÀ ÇÂIñîÊ | | RY: ÈÄÅÍÒÈÔÈÈÀÖÈß |

’Kðèñòà | RY: ÈÄÅÍÒÈÔÈÈÀÖÈß |

çäðââñòâóéòâ? | RIE: ÍÐÈÅÅÒÑÒÅÈÅ |

(.)

H: çäðââñòâóéòâ, | RIJ: ÍÐÈÅÅÒÑÒÅÈÅ |

=y `õiòåë áû ((íacâàíèå ôëðiû)) (.) < äèñ`iåò÷åðà > èëè `êàëîé-íèáóäü `íi- íiâð. | DIE: ÄÈÐÅÈÒÈÅ |

(0.8)

V: i ñèíóòî÷éó. | DIJ: ÍÒÑÐÍxÉÀ |

(5.5)

V: òåëåôíñ iá` ñëóæèâàíèý [òî åñòü,] äà? | KYE: ÍÁÙÈÉ ÅIIÐÍÑ | VTE: ÓÒÍxÍÅÍÈÅ ÓÑËÍÂÈÉ ÍÒÅÅÒÀ |

H: [äà.] | KYJ: ÄÀ | VTJ: ÓÒÍxÍÅÍÈÅ ÓÑËÍÂÈÉ ÍÒÅÅÒÀ |

V: ï íåñäòåðæääííùì äàííùì ïíåð < ÷åòûðå ÷åòûðå ñåìü? (0.5) äåâýöü åññåìü? (.) [ïäèí] ïíëü. > | DIJ: ÈÍÔÐÌÀÖÈß |

H: [.hh] | YA: ÈÍÅ |

(.)

H: ñiiñèáí. | RIE: ÑIIÑÈÁÍ |

V: ïæàëóéñòà? | RIJ: ïIÆÀËÓÉÑÒÀ |

2.3. Êíííóíèéàòèâíúå ñòðàòåãèè

lú èññõtæì èç iíjyöèý êíííóíèéàòèâííé ñòðàòåãèè è êííñòðóêòèâííé iíäåéè äèàëíà (Constructive Dialogue Model, CDM) [Jokinen 0-àñòíèéíí áæàëíà äëý iíñòðíåíèý ñëåäóþùåé ðåïëèéè êâé ðåàéëè íà iðåäûäóùóþ ðåïëèéó íàðòíåðà. Êíííóíèéàòèâíúå ñòðàòåãèè íáñáùáíííí óðíñáíá.

Äëý iíðåäåéäíèý êíííóíèéàòèâíúô ñòðàòåãèé â CDM èññëüçóþòñý ÷åòûðå êííòåêñòóàëüíûô ôàéòíðà:

– íæèääååññòü ðåïëèéè;

– ñâýçü ðåïëèéè ñ òåííé;

– äñòèäíóòñòü öåëåé åíâíðýùååí;

– èíèöèàòèâà åíâíðýùååí èëé íàðòíåðà.

Âñå êííòåêñòóàëüíûå ôàéòíðû áèíàðíû, ñíñòååòñòååíííí, êíííóíèéàòèâíúô ñòðàòåãèé $2^*4=16$.

Íàèáíëåå óää÷ííé ñòðóêòóðíé iðåäñòàåëåíèý êíííóíèéàòèâíúô öåëåé ýâëÿðòñý íàäàçèí öåëåé (stock) è íðèíöèí «íññëåäíèí áîøåé, íññëåäíàòååëüíí íäíá íäá äðóäíé. Íí íàðå åñòèæåíèý ååðöíéò öåëåé ííé óääëýþòñý; áñå öåëè äèàëíà ñ÷èòàþòñý åñòèäíóòûé,

Êíííóíèéàòèâíúå ñòðàòåãèè ðàçìå÷åíû á 20 åñòåñòååííûô è 20 ñèíóëýöèííûô äèàëíàð.

Iðèìâð 2 (ñèlóëyöèíííûé äèàëíâ â iàðåââíâá íà ðóññéèé ýçûê).

Dâëëëëà

Nòðàòåãëÿ

Iàãàçèí

Êëëåíò:

Èàë äâñàòòü äî ïýðíó èç Òàðòó äî ïëóäíý?

Íæëääâái.-ñâýçí.-äîñòèãí.-äîâîðýù.-íà÷àëí/êííâö

–

Êññüþòåð:

Àâòîáóñ âûõíäèò à 5 óòðà.

Àâòîáóñ âûõíäèò à 8 óòðà.

Íæèäàåì.-ñâýçí.-íáäíñòèäí.-ãîâîðýù.-íðîäíëæåíèå òåìû

Òàðòó-ÿðíó

Èíòåðåñóåò ëè âàñ âðåìÿ íðèáûòèÿ? ...

Íáíæèäàåì.-ñâýçí.-äíñòèäí.-ãîâîðýù.-íâûé äèàëíâ

Òàðòó-ÿðíó

Đàciā-áííûé êîðiōñ èññiēúçóâòñý â òâiðâòè÷âñêèõ è iðeêëàäíûõ èññéâäiâàíèýõ (íàïð., [Hennoste et al. 2005]) è â ðàáîòâ iàä à [Fishel 2004]), à òâêæâ â iññòðiáíèè iâðâûõ äèàëiâûõ ñèñòâi à yñòíñêi ýçûêå. Ýòè ñèñòâiû, Òðàíññðòíûé è Òáàòðàëüíûé À iðeìâiáíèâi çíáíèé i ñòðóêòóðå èíôiðiàöèííûõ äèàëiâi è iññäèëiâi.

Ñièñîê ëèòåðàòóðû

Mark Fishel 2005. Dialogue Act Recognition in Estonian Dialogues Using Artificial Neural Networks. In: Proceedings of the Second Baltic Conference on Human Language Technologies, Tallinn, 4–5 April 2005, 249–254.

Olga Gerassimenko, Tiit Hennoste, Mare Koit, Andriela Rääbis, Krista Strandson, Maret Valdisoo, Evely Vutt. Annotated Dialogue Corpus as a Language Resource: An Experience of Building the Estonian Dialogue Corpus. The First Baltic Conference “Human Language Technologies. The Baltic Perspective”. Commission of the Official Language at the Chancellery of the President of Latvia, Riga, 2004, 150–155.

Tiit Hennoste, Olga Gerassimenko, Riina Kasterpalu, Mare Koit, Andriela Rääbis, Krista Strandson, Maret Valdisoo. Questions in Estonian Information Dialogues: Form and Functions. Text, Speech and Dialogue. 6th International Conference TSD 2005. Springer, 2005, 420–427.

Gail Jefferson 2004. Glossary of transcript symbols with an introduction. In Lerner, G.H. (Ed). Conversation Analysis: Studies from the first generation. Amsterdam/Philadelphia: John Benjamins, 13–31.

Kristiina Jokinen. 1995. Rationality in Constructive Dialogue Management; URL: <http://cl.aist-nara.ac.jp/lab/papers/kris/aaai.ps> (used 20.05.2006).

Áëàäiâðiñòè

Đàáîòó iññäâðæèâàåò Ýñòíñêèé iàó÷íûé ôííä (âðàíò 5685).

Summary

The Estonian Dialogue Corpus is collected with the aim of developing the dialogue system using the natural language. The spoken dialogues (884 dialogues, 155000 running words) are used to study the rules and norms of the human-human communication; the corpus also includes human-computer dialogues (21, 2500 running words) collected by the Wizard of Oz method used to study the role behaviour of the users and information provider. The presentation considers the means and levels of transcription and annotation dialogues and also the application of the corpus.

Iðeëiæåíèå

1. Öðàíñêðèëöéííúå çíàéè

ñiiää èíòííàöèè

ïíëóñiiää èíòííàöèè

ïíäúåì èíòííàöèè ?

êîðîòêàÿ iàóçà (iàéñ. 0.2 ñ.) (.)

äëèíà iàóçû â ñâêóíäàõ (2.0)

iàéíæåíèå [text]

ñëëÿíèå íåçàâèñèíûõ åäëíèö =

ðàñòÿíóòûé çâóê ::

óääðííà ñëíâí

íðåðâàííà ñëíâí do-

âäîõ .hhh

óáúñòðåíèå òåííà > text <

çàìåäëåíèå òáiíà < text >

ñíííèòåëüíûé ìòðåçíê {text}

íåðàçáíð÷èåûé ìòðåçíê {---}

2. Òèïíèíæý ðå÷åûõ àêòíâ

I Åêòû, ñíñòàâëýþùèå ñìåæíûå íaðû

1.1 Åêòû óïðàâëåíèý äèàëíâí

1.1.1 Kíílóíèêàöëý

1.1.1.1 Ðèòóàëü (RIE RIJ)

1.1.1.2 Niñáà òåìû

1.1.2 Ðàçõåøåíèå ïõláëåì

1.1.2.1 Èñïðàâëåíèå, èíèöèèõlâàííà ïàðòíåðî

1.1.2.2 Íðlâåðêà êííòàêòà

1.1.2.3 Óòí÷íáíèå óñëíâèé ìòâåòà (VTE VTJ)

1.2 Èíòíðìàöèííûå àêòû

1.2.1 Äèõåêòèåû (DIE DIJ)

1.2.2 Âññðñû (KYE KYJ)

KYE: íáùèé âññðñ

KYE: íáùèé âññðñ, íæèäàþùèé ðàçâåðíóòíâí ìòâåòà

KYE: àëüòåðíàòèâíûé âïïðîñ

KYE: ñïåöèàëüíûé âïïðîñ

KYE: èííå

KYJ: äà

KYJ: íåò

KYJ: ñîãëàñííå íåò

KYJ: èííé îòâåò íà íáùèé âïïðîñ

KYJ: àëüòåðíàòèâà: íäíà

KYJ: àëüòåðíàòèâà: íáå

KYJ: àëüòåðíàòèâà: òðåòèé âûáîð

KYJ: àëüòåðíàòèâà: îòðèöàíèå

KYJ: àëüòåðíàòèâà: èííå

KYJ: ðàçâåðíóòûé îòâåò

KYJ: îòñóòñòâèå èíôîðìàöèè

KYJ: îòêàç

KYJ: èííå

1.2.3 Ííåíèå

II Íäèí÷íûå àêòû

2.1 Àêòû óïðàâëáíèÿ äèàëíâíí

2.1.1 Èííóíèêàöèÿ

2.1.1.1 Đèòóàëû (RY)

2.1.1.2 Íáðàòíàÿ ñâýçü

2.1.2 Ðàçðåøåíèå ïðíáæå

2.1.2.1 Èñïðàâëåíèå

2.2 Èíôîðìàöèííûå àêòû

2.2.1 Íñííâíûå àêòû (YA)

2.2.2 Äññëíåíèÿ íñííâíûõ àêòîâ (ííyñíåíèå, óòî÷íåíèå)

Íðèìå÷àíèå: Íäðíáíí ðàññëñàíà ãðóíà ãíñðîñâ è ìòâåòîâ —; íàèáíëå ÷àñòòíûõ ðå÷åâûõ àêòîâ; ýñòííñëàÿ àááðåâèòóðà,

Summary

The Estonian Dialogue Corpus is collected with the aim of developing the dialogue system using the natural language. The spoken dialogues (884 dialogues, 155000 running words) are used to study the rules and norms of the human-human communication; the corpus also includes human-computer dialogues (21, 2500 running words) collected by the Wizard of Oz method used to study the role behaviour of the users and information provider. The presentation considers the means and levels of transcription and annotation dialogues and also the application of the corpus.