
Eesti Keeleressursside Keskus

Kadri Vider (kadri.vider@ut.ee)

Eesti Keeleressursside Keskus
www.keeleressursid.ee



Eesti Keeleressursside Keskus (EKRK)

- Eesti keeleressursside keskus on humanitaaride teadustaristu – erinevates uurimisasutustes paiknevate, veebist ligipääsetavate andmehoidlate võrgustik, mis võimaldab autentimise teel juurdepääsu mitmel erineval tasemel kasutajatele.
- Lisaks olemasolevate ja uute, loodavate keeleressursside kogumisele ja arhiveerimisele käivitatakse süsteem olemasolevate keeleressursside tutvustamiseks ja potentsiaalsete kasutajate koolitamiseks.

EKRRK missioon

□ milleks?

- Luua taristu, mis võimaldaks kõigile uurijatele **keeleressursside ja -tehnoloogiate kättesaadavuse**
- Keelest sõltumatuid vahendeid on võimalus kasutada ja jagada
- Keelest sõltuvaid vahendeid on võimalik üle kanda

□ kuidas?

- Ühendades eksisteerivad digitaalsed arhiivid ja tagades nende kättesaadavuse veebi kaudu
- Pakkudes keeletehnoloogia vahendeid kui veebiteenust, mis kasutab arhiveeritud andmeid
- Meta-andmeid kättesaadavaks tehes

EKRRK kui konsortsium

- konsortsiumleping allkirjastati 2.12.2011



TARTU ÜLIKOOL



TTÜ KÜBERNEETIKA INSTITUUT
Institute of Cybernetics at TUT



Eesti Keele Instituut

Mõisteid

□ Keeleressurss

- andmestik keele kohta või vahend/abinõu sellise andmestiku esitamiseks või töötlemiseks
- näiteks tekstid, salvestised, keelekirjeldused ja – annotatsioonid, välitööde materjalid, tarkvara, protokollid, andmemudelid, veebiarhiivid ja -indeksid
- digitaalne või mitte, avaldatud või käsikirjas

□ Meta-andmestik

- Kirjeldav struktureeritud info ressursi kohta, näiteks info teavikute kohta kataloogikaardil raamatukogus, digimaailmas sageli kasutajale peidetud märgenditega info digiobjekti kohta.
- Sisaldab kokkulepitud elementide komplekti (näiteks DCMI, IMDI, IsoCAT), mis võimaldab ka otsingut filtreerida

Digitaalsed andmed

The image shows a digital library interface. On the left, a sidebar titled "Pages" displays a grid of 10 document thumbnails, numbered 1 through 10. Thumbnail 1 is highlighted with a red border. The main area on the right shows a large preview of the selected document (Thumbnail 1). The preview has a white header with the logo "eod | books2ebooks.eu". Below the header is a blue horizontal bar. The main content area is light gray and contains the text "HALDUR ÕIM" and "Inimene, keel ja arvuti ehk kompuuterlingvistika".

Mis on keeleressurs?

www.filosoft.ee/html_morf_et/html_morf.cgi

Most Visited Täna ülikoolis - ut.ee Introducing

olnu+d // _S_ pl n, //

infootsingumeetodite
info_otsingu_meetod+te // _S_ pl g, //

analüüs
analüüs+0 // _S_ sg n, //

Kuid
kuid+0 // _J_ //
kuu+id // _S_ pl p, //

raamatukogude
raamatu_kogu+de // _S_ pl g, //

kui
kui+0 // _D_ //
kui+0 // _J_ //

näite
näi+te // _V_ te, //
näide+0 // _S_ sg g, //
näit+e // _S_ pl p, //

varal
vara+l // _S_ sg ad, //
varal+0 // _K_ //

6im_inimekeelarvuti.pdf - Adobe Reader

File Edit View Document Tools Window Help

33 / 149 100%

infootsing

2. Tekstide otsing ja automatiseeritud infootsing

Erinevalt inimese ajast ei unusta arvuti vastava signaali saamisel kõiki jälle välja laduda. Seepärast ongi ehk kõige kasulikum arvutile üldse mitte midagi ütelda.

C. Ford, «Mõtlemise õpetus»

2.1. Automatiseeritud infootsisüsteemid – probleemid ja ehitus

Vajamineva informatsiooni ülesotsimine kui omaette probleem sündis ilmselt ühel ajal esimeste arhiivide ja raamatukogude tekkega, seega 4–5 aastatuhandet tagasi. Juba kolme ja poole tuhande aasta eest valitsenud vaarao Ramses II raamatukogu sisaldanud 20 000 papüürusrulli, nii et vajaliku informatsiooni leidmine ei saanud tallegi olla triviaalne ülesanne.

Möödaläinud aastatuhandete jooksul on välja töötatud küllalt tõhusad meetodid ja vahendid valialike raamatute alakiriade ja ar-

tarbijate üha pakilisemaid ja spetsiifilisemaid vajadusi. Meie eesmärgiks ei ole muidugi raamatukogude töötamis põhimõtete kirjeldamine ega ka neis seni kasutusel olnud infootsingumeetodite analüüs. Kuid raamatukogude kui näite varal jõuame kergemini meid huvitavate probleemide juurde.

Kui meid huvitab mingi kindla küsimuse kohta leiduv kirjandus, siis peamine allikas selle väljasegitamisel raamatukogus on süstemaatiline kataloog. Selles on kogu raamatukogus leiduv kirjandus klassifitseeritud sisu järgi kindlatesse liikidesse, mis on hierarhilise ehitusega, s. t. jagunevad järjest kitsamateks alaliikideks. Tänapäeval on enamikus maailma maa-des kasutusel ühtne teadus- ja

TEKST: Meie eesmärgiks ei ole muidugi raamatukogude töötamis põhimõtete kirjeldamine ega ka neis seni kasutusel olnud infootsingumeetodite analüüs.



Humanitaarteadlaste andmekogud

- Tänapäevaks on digitaliseeritud tohutu maht humanitaarteaduste uurimisandmeid, enamik nendest on keelepõhised
- Paljud sellised arhiivid kasutavad erinevaid standardeid, sõltuvalt uurimise eesmärgist on andmed erineva detailsuse või struktuuriga
- Ka andmetele ligipääs on korraldatud eri viisidel
- Humanitaarteadlased sageli ei tea
 - mis on keeleressursid (KR)
 - kas ja kuidas KR neid võiks aidata

EKRRK Euroopa plaanis

- **CLARIN** - www.clarin.eu
Common Language Resources and Technology Infrastructure
 - Tugeva kasutajatoega, teadlastele orienteeritud võrgustik
 - Koosneb eri tüüpi keskustest, mis ühendavad kasutajaid, ressursse ja tugiteenuseid
 - EKRRK on Eesti CLARINI keskus

- **META** - www.cs.ut.ee/metanord/
 - Ühtne repositooriumite võrgustik Euroopa mitmekeelsuse toeks
 - Võrgustiku sõlmed vahendavad ressursse ühtses metaandmete vormis ja ühtsetel tingimustel

- **DASISH** - www.dasish.eu
 - Sotsiaal- ja humanitaarteaduste andmeteenuste taristu
 - Võrgustikus CLARIN, DARIAH, ESS, CESSDA, SHARE

CLARIN LRT inventory: Estonian

[VLO Home](#) >> Faceted Browser Resources

Estonian

search

COLLECTION : [all](#) > CLARIN LRT inventory

LANGUAGE

[Estonian](#) (29)
[English](#) (7)
[Hungarian](#) (6)
[Bulgarian](#) (5)
[Czech](#) (5)
[Romanian; Moldavian; Moldovan](#) (5)
[Danish](#) (4)
[Dutch](#) (4)
[Finnish](#) (4)
[French](#) (4)
[more...](#)

CONTINENT

COUNTRY

[Estonia](#) (22)
[Germany](#) (2)
[Latvia](#) (2)
[Bulgaria](#) (1)
[Canada](#) (1)
[Croatia](#) (1)
[Czech Republic](#) (1)
[Hungary](#) (1)
[Italy](#) (1)
[Lithuania](#) (1)
[more...](#)

ORGANISATION

[University of Tartu](#) (8)
[Institute of Cybernetics at Tallinn University of Technology](#) (2)
[Department of Linguistics and Nordic Studies, University of Oslo](#) (1)
[Joint Research Centre of the EU](#) (1)
[Netherlands Institute of Advanced Study \(NIAS\), Collegium Budapest](#) (1)
[Tilde](#) (1)
[Tilde as coordinator, Eurotermbank consortium](#) (1)

GENRE

SUBJECT

RESOURCETYPE

[Written Corpus](#) (10)
[Spoken Corpus](#) (5)
[Lexicon / Knowledge Source](#) (4)
[Treebank](#) (3)
[Terminological Resource](#) (2)
[Aligned Corpus](#) (1)
[Grammar](#) (1)
[Web Service](#) (1)

Showing 1 to 10 of 29

<< < 1 **2** 3 > >>

name	description
Arbores	149 sentences, VISL tagset
BABEL Estonian Database	The database consists of three sets: - Many Talker Set: 30 males, 30 females; each to read 50...
Constraint Grammar of Estonian	general written, Constraint Grammar
Corpus of Old Written Estonian	Corpus of texts written fully or partly in Estonian, from 13.-19. century; 1,5 million words
Corpus of Present-day Written Estonian	written general; 95 mio words; TEI/SGML
Corpus of Spoken Estonian	spoken general; 1 mio words; local tagset
Corpus of Written Estonian	4.4 mio words; TEI/SGML
Corpus query for Estonian corpora	Web application for querying the automatically morphologically disambiguated Mixed corpus of...
Corpus with Disambiguated Word Senses	100000 words, word senses based on TEKsaurus (Estonian Wordnet)
Database of Estonian Multi-word Verbs	17 500 entries

META-SHARE repositooriumide võrgustik

- META-NETi võrgustiku (*A Network of Excellence forging the Multilingual Europe Technology Alliance*, <http://www.meta-net.eu/>) hajustaristu, mille metadata registri (www.meta-net.eu/meta-share) ja repositooriumi andmetele pääseb vabalt ligi võrgustikku kuuluvate sõlmede (node) kaudu.
- Tartu Ülikool osales META-NETi projektis META-NORD, mille üks ülesandeid panna püsti oma võrgustikusõlm metaandmete registri ja repositooriumiga: metashare.ut.ee

META-SHARE andmestik

- Registri andmeteks on repositooriumis säilitatava keeleressursi meta-andmed:
 - Esindatud keeled, kas üks või mitmekeelne ressurss, paralleelne vm struktuur mitmekeelsetel
 - Ressursi tüüp, meediumi ja esitusviisi tüüp (tekst/audio..., kirjalik/suuline keel...), formaat
 - Kättesaadavus
 - Kasutuslitsentsi tüüp ja kasutajaskonna piirangud (k.a inim- või masinkasutatavus)
 - Võimalikud KT rakendusvaldkonnad



Keywords:

[Return to Browse page](#)

Search

Filter by:

Language:

- ✦ Estonian (16)
- ✦ English (5)
- ✦ German (5)
- ✦ Latvian (5)

[more](#)

Resource Type:

- ✦ corpus (14)
- ✦ lexicalConceptualResource (12)

Media Type:

- ✦ text (21)
- ✦ audio (6)

Availability:

- ✦ available-restrictedUse (23)
- ✦ available-unrestrictedUse (3)

Licence:

- ✦ ELRA_END_USER (13)
- ✦ ELRA_VAR (13)
- ✦ ELRA_EVALUATION (7)
- ✦ proprietary (4)

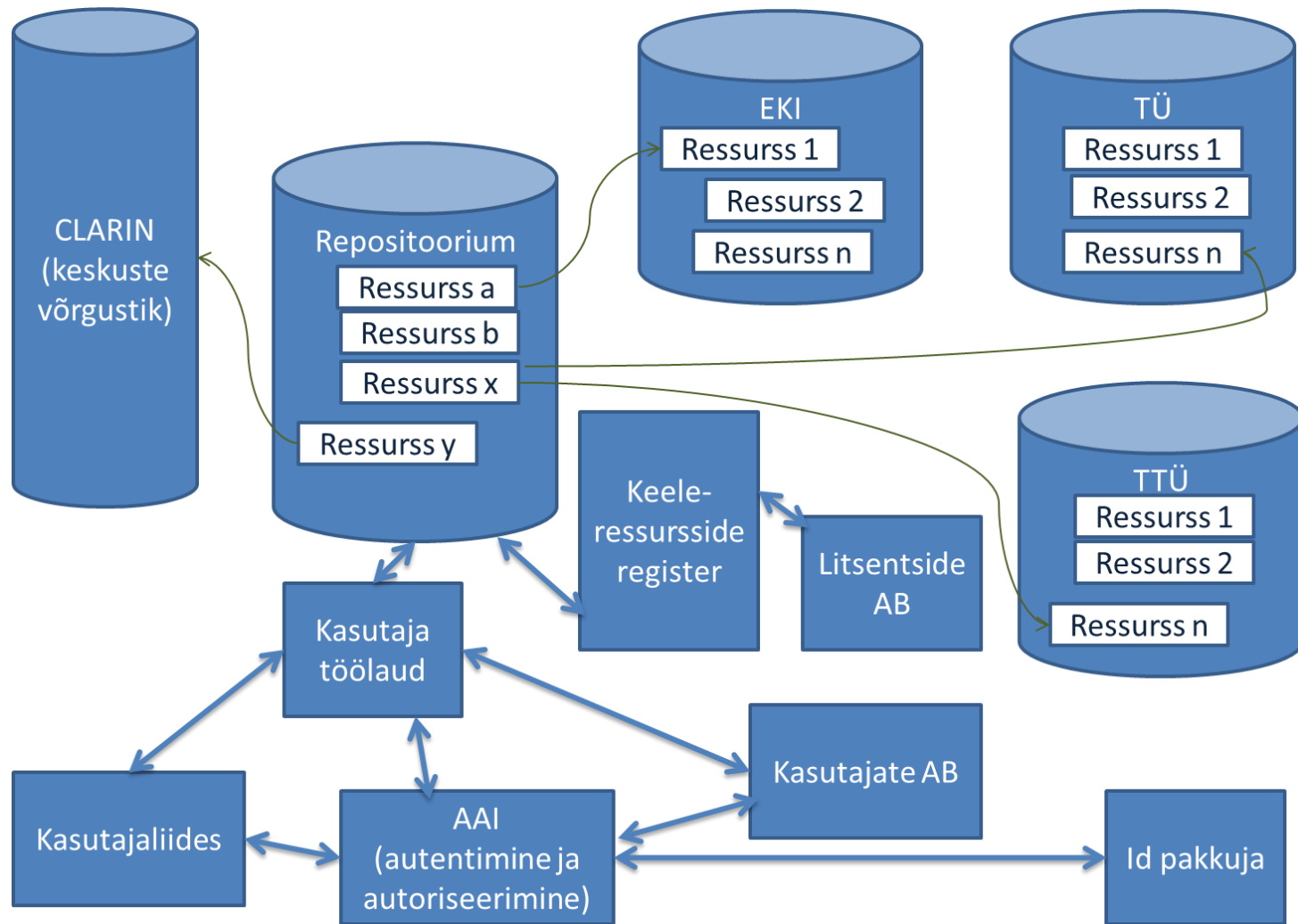
[more](#)

26 Language Resources (Page 1 of 2)

[« Previous](#) | [Next »](#)

Resource Name	Resource Type	Media Type	Language
ACCURAT Bilingual Comparable Corpora for under-resourced languages	corpus	text	Croatian, English, Estonian, German, Italian, Latvian, Modern Greek (1453-), Romanian, Slovenian
BABEL Bulgarian Database	corpus	audio	Bulgarian
BABEL Estonian Database	corpus	audio	Estonian
BABEL Hungarian Database	corpus	audio	Hungarian
BABEL Polish database	corpus	audio	Polish
BABEL Romanian database	corpus	audio	Romanian
Bulgarian X language Parallel Corpus Bul-X-Cor	corpus	text	Albanian, Bosnian, Bulgarian, Croatian, Czech, Danish, Dutch, English, Estonian, Finish, Galician,

Keskuse komponendid



Keskuse pakutavad teenused

- Keeleressursside arhiveerimine ja haldamine
- Keeleressursside kogumine ja hindamine
- Ligipääs ja kasutajate koolitamine

Avatud nii keeleressursside pakkujatele, arendajatele kui ka keeleressursside kasutajatele, kes nõustuvad kasutuskorra ja litsentsitingimustega, kuid eelistatud on teadus-arendusasutuste kasutajad ja partnerid CLARINi liikmete seas.

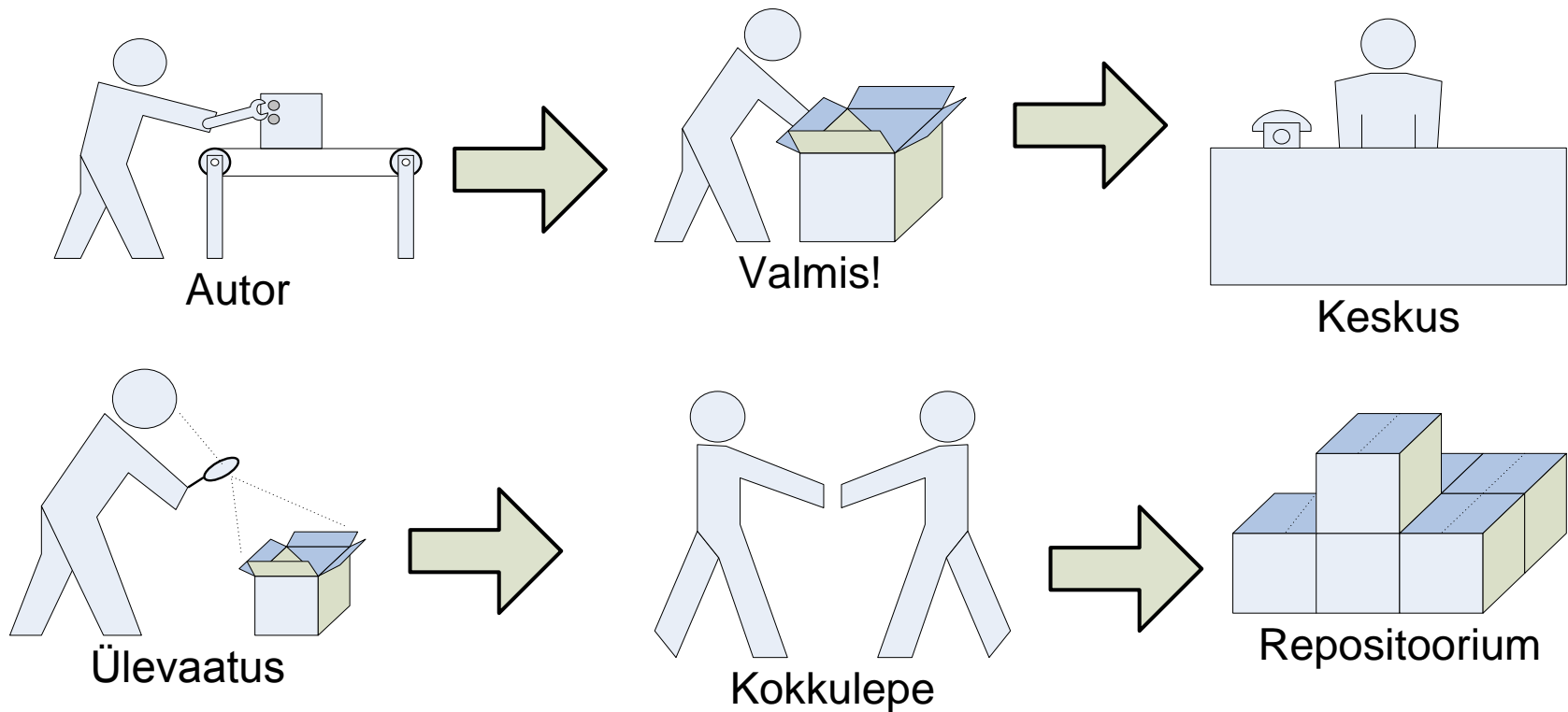
Arhiveerimine

- ressursside metaandmete (ehk tekstilise kirjelduse) säilitamine registris
- ressursside koopiate säilitamine repositooriumis

Ressursside kogumine

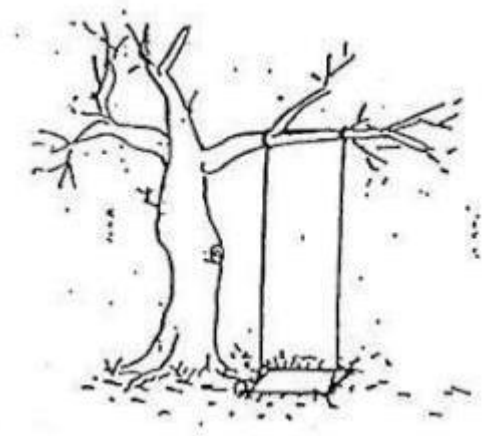
- Keeleressursside keskuse partnerid
 - Koondada olemasolevad tööruhmade ressursid
 - Hoida kättesaadaval ressursside vanemaid versioone
 - Ressursside ja tarkvara kasutamine üle veebi
- EKKTT programmi raames loodud ressursid
 - Loodud ressursid säilitatakse ühtses kohas
 - Võimaldatakse ressurssidele juurdepääs vastavalt kasutustingimustele
 - Võimalus ressursse edasi arendada
- Kõik teised soovijad

Kuidas ressursss meile jõuab? 😊



Ligipääs

- Ressurssidel 3 tüüpi kasutuslitsentse
 - Vaba kasutus kõigile (näiteks Creative Commons)
 - Kasutamiseks teadustöö eesmärkidel (ACA)
 - Kasutamiseks eritingimustel (mitte-kommerts või isikuandmetega seotud)
- Kasutajate võimalused sõltuvalt kuuluvusest
 - Laialdasimad konsortsiumipartneritel
 - CLARINi partnerid jt teaduskasutajad
 - Avalikkus



Ligipääs

- Luuakse avalik veebipõhine ligipääs
 - www.keeleressursid.ee
- Luuakse kasutajagrupid, määratakse kasutusõigused.
- Eelisolukorras ligipääsu võimaldamisel arhiveerimisteenusele on CLARIN-ERIC liikmed.
- **Liidestus CLARINiga**
 - Regulaarne andmevahetus
 - SSO* autentimine, ligipääs rahvusvahelisse võrgustikku

Koostöö mäluasutustega (1)

- Digiteerimise ja kopeerimise alane koostöö (ka juriidilisest aspektist)
 - tekstilise materjali osas huvitab meid ainult OCR-tud ehk tähtsustatud materjal
 - helimaterjali analüüsiks ja massiliseks sisuotsinguks võimalik vastastikku kasulik koostöö KübI kõnetehnoloogia spetsialistidega
- Arhiveerimise ja pikaajalise säilitamise alane koostöö
 - varukoopiate deponeerimine
 - PID-süsteemi sünkroonimine

Koostöö mäluasutustega (2)

- Sisu kasutamise autoriõiguse ja litsentsimise teemad
- Sisule ligipääsu teema, sealhulgas koostöö kasutajate autentimise ühise süsteemi alal
- Sisu semantilise annoteerimise teemad otsisüsteemide jaoks – see on ka keeletehnoloogiline ülesanne

Täna tähelepanu eest!

