

AUTOMATIKUS SZÓTÁRÉPÍTÉSI MÓDSZEREK FINNUGOR NYELVEKRE

Ferenczi Zsanett

2020. június 25.

PPKE BTK Nyelvtudományi Doktori Iskola

Nyelvtechnológia Műhely

Témavezető: Simon Eszter

1. A kutatás bemutatása
2. Lexikográfiai adatbázis
3. Szócikkek felépítése
4. További feladatok

A kutatás bemutatása

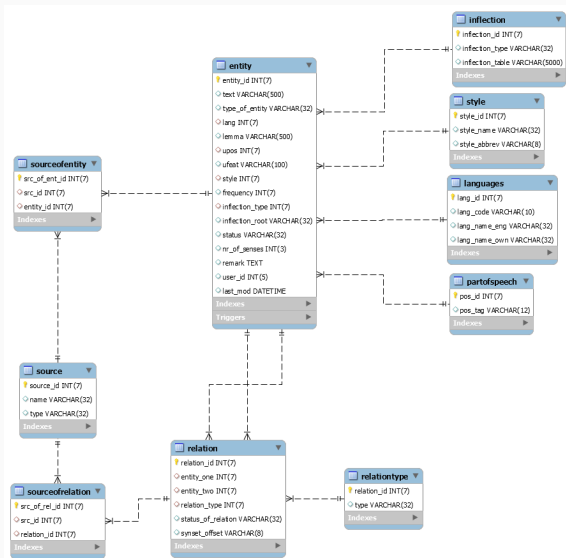
- *"Automatikus szótárépítési módszerek finnugor nyelvekre"*
- szótárépítési módszerek összehasonlítása
- finnugor nyelvek többnyelvű online szótára
- nyelvtanuló alkalmazás (CALL)
- létező nyelvtechnológiai eszközök közötti kompatibilitás biztosítása, új eszközök létrehozása

- **új eszközök:**
 - finn és magyar Wiktionaryk feldolgozása
 - finn és magyar WordNetek feldolgozása, összekötése
 - tremjugani és vaszjugani hanti > betűrendbe sorolás
- **szótár:**
 - szópárok kinyerése különböző módszerekkel
 - nyelvfüggetlen lexikográfiai adatbázis felépítése
 - lemmák morfológiai elemzése, szófaji taggelés
 - definíciók, inflexiós tőtípusok, példamondatok kinyerése

Lexikográfiai adatbázis

- cél: **többszavú szótár** adatainak elraktározása, összegyűjtött információk **redundanciamentes** tárolása
- minden nyelvi információt azonos módon kezelünk → **entitás**
- az entitások típusai: *lemma, szóalak, többszavas kifejezés, definíció, mondat, stb.*
- a jelentést az egyes entitások közötti relációkkal tudjuk megragadni
- a relációk közt megtalálhatók szemantikai relációk (mint *szinonima, antonima*), illetve különböző lexikográfiai relációk (pl. *fordítása, definíciója, példamondata, stb.*)
- az entitáspárok és entitások forrását eltároljuk

Az adatbázis



The image shows a screenshot of a database table definition for a table named 'entity'. The table has the following columns and data types:

- entity_id INT(7) (Primary key, indicated by a yellow lightning bolt icon)
- text VARCHAR(500)
- type_of_entity VARCHAR(32)
- lang INT(7)
- lemma VARCHAR(500)
- upos INT(7)
- ufeat VARCHAR(100)
- style INT(7)
- frequency INT(7)
- inflection_type INT(7)
- inflection_root VARCHAR(32)
- status VARCHAR(32)
- nr_of_senses INT(3)
- remark TEXT
- user_id INT(5)
- last_mod DATETIME

Below the table definition, there is a section labeled 'Indexes' with a right-pointing arrow. At the bottom of the screenshot, there are two vertical lines with horizontal bars at the top, representing a primary key constraint.

Szócikkek felépítése

apu  **noun**

isä, iskä

(Apa megszólítása gyermeke részéről)

Apu szereti anyut.

Isä rakastaa äitiä.

szinonimák: **apuci, apuka**

apu **+** **noun**

segítség

Apu tuli viime hetkellä.

A segítség az utolsó pillanatban érkezett.

ég  **noun**

taivas

etimológia: *Ősi finnugor örökségünk: zürjén szünöd ('levegő'), finn sää ('időjárás').*

szinonimák: **égbolt**, **menny**

összetételek: **éghajlat**, **égtáj**

ég  **verb**

palaa

(Hő és fény kibocsátása közben pusztul, fogy.)

ragozás:

E/1 égek

E/2 égysz

E/3 ég

T/1 égünk

T/2 égtek

T/3 égnek

...

szinonimák: **lángol**, **parázslík**

daru  **noun**

(madár) **kurki**

(szürke, gólyaszerű gázlómadár)

etimológia: *ogul tárov, osztják táreh, zürjén, votják turi* ('daru').

ragozás: ... darvak ...

(emelőgép) **nosturi**

ragozás: ... daruk ...

finn és magyar szópárok

- Wiktionary - **Wikt2dict**
 - Ács et al. (2013) Wiktionaryt feldolgozó eszköze
 - általában lemmák, előfordulnak nem alapalakok is
 - két módon: *extract*, *triangulate*
 - finn, magyar és angol Wiktionary verziók
 - összesen **16 685 egyedi szópár**
- finn és magyar Wiktionary
 - saját fejlesztésű, Wiktionary dumptot feldolgozó eszköz
 - szópárokat, szófaji címkéket és tőtípust nyer ki
 - összesen **8 762 egyedi szópár**

Összegyűjtött szópárok

- WordNet
 - FinnWordNet (Lindén és Carlson (2010))
 - HuWN (Miháltz et al. (2008))
 - **99 780 szópár** (szófaji címkével)
 - 25 662 synsetpár

poikanen, kissanpentu cica, kiscica, kölyökmacska, kismacska

- OpusCorpus
 - párhuzamos korpusz
 - szavak nincsenek lemmatizálva
 - **832 696 szópár**

kissanpennuilla macskakölyöknek

kissanpennut kiscicák

kissanpentu kiscica

- **Wikt2dict** eszköz és **OpusCorpus** nem ad meg szófaji címkéket
- morfológiai elemzőket használva kell szófajt kinyerni
- egységes UD kimenetettel
- **omorfi** - finn szavakhoz
- **emMorph** - magyar szavakhoz, **emMorph2UD** eszköz
- minden elemzést megtartva

- magyar Wiktionaryról magyar szavakhoz
 - 27 154 egyedi szóhoz 31 469 definíció

pad Fából vagy kőből készült hosszú ülőke.

- finn Wiktionaryról finn szavakhoz
 - 102 638 egyedi szóhoz 127 786 definíció
- magyar WordNetről
 - 35 011 egyedi szóhoz 42 356 definíció

dolog Független, önmagában álló entitás.

dolog Valamilyen cselekedet, tett, magatartás.

dolog Elvégzendő feladat, munka.

- példamondatok a WordNetről:
 - magyar: **22 562** synsethez
- paradigmák
 - finn Wiktionaryről: **24 121** finn szóhoz
 - magyar Wiktionaryről: **4 704** magyar szóhoz
 - tőtípusok sablonjai:

poliszémia szócikk esetén:

```
{{hu-fn-k1|poliszémi|á}}{{hu-birt-tok|poliszémiá}}
```

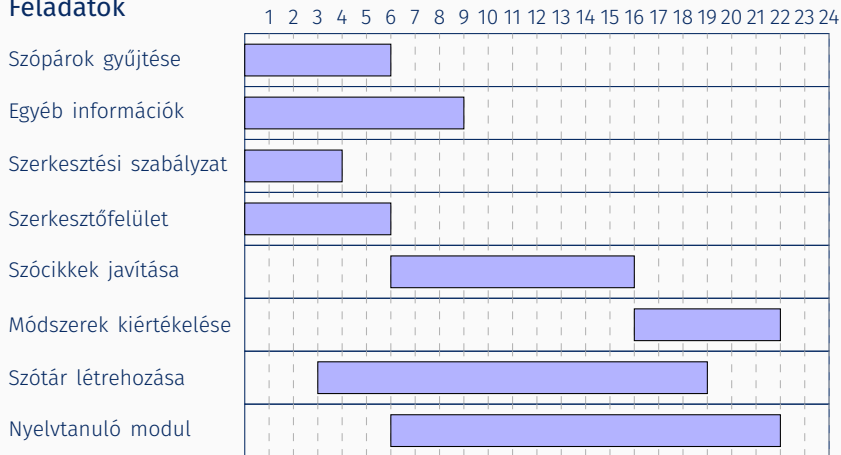
- a tőtipusokhoz tartozó inflexiós táblák kinyerése

<i>hu-fn-k1</i> ragozása [becsuk]		
eset/szám	egyes szám	többes szám
alanyeset	hu-fn-k1	{{1}}>{{2}}k
tárgyeset	{{1}}>{{2}}t	{{1}}>{{2}}kat
részes eset	{{1}}>{{2}}nak	{{1}}>{{2}}knak
-val/-vel	{{1}}>{{2}}val	{{1}}>{{2}}kkal
-ért	{{1}}>{{2}}ért	{{1}}>{{2}}kért
-vá/-vé	{{1}}>{{2}}vá	{{1}}>{{2}}kká
-ig	{{1}}>{{2}}ig	{{1}}>{{2}}kig
-ként	hu-fn-k1ként	{{1}}>{{2}}kként
-ul/-ül	-	-
-ban/-ben	{{1}}>{{2}}ban	{{1}}>{{2}}kban
-on/-en/-ön	{{1}}>{{2}}n	{{1}}>{{2}}kon
-nál/-nél	{{1}}>{{2}}nál	{{1}}>{{2}}knál
-ba/-be	{{1}}>{{2}}ba	{{1}}>{{2}}kba
-ra/-re	{{1}}>{{2}}ra	{{1}}>{{2}}kra
-hoz/-hez/-höz	{{1}}>{{2}}hoz	{{1}}>{{2}}khoz
-ból/-ből	{{1}}>{{2}}ból	{{1}}>{{2}}kból
-ról/-ről	{{1}}>{{2}}ról	{{1}}>{{2}}król
-től/-től	{{1}}>{{2}}től	{{1}}>{{2}}ktől

További feladatok

Kutatással kapcsolatos teendők

Feladatok



- feleletválasztós kérdések
 - pl. _____ *on iso kissa.*
 - a)** eläin **b)** hiiri **c)** koira **d)** tiikeri
 - definíciókból, példamondatokból
 - a célszóhoz hasonló választási lehetőségekkel
- mondatba illő szóalak megadása
 - pl. *Mitä nähtävyyksiä _____ (kaupunki) on?*
 - definíciókból, példamondatokból
- szókincsfejlesztés
 - forrás- és célnyelvi szavak összepárosítása
 - szinonimák keresése

- Recenzió magyar nyelvű könyvről
- Magyar nyelvű lektorált kiadvány szerkesztése
- Tudománynépszerűsítő közlemény
- Egyéb elismert tudományos vagy tudománynépszerűsítő tevékenység
- Oktatási modul
- Saját kutatáson alapuló tudományos publikációk
- Konferenciaelőadások (AlkNyelvDok, MSZNY, IWCLUL, NLP4CALL)
- Posztterek

Judit Ács, Katalin Pajkossy, and Andras Kornai. Building basic vocabulary across 40 languages. In *Proceedings of the Sixth Workshop on Building and Using Comparable Corpora*, pp. 52–58, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL

<http://www.aclweb.org/anthology/W13-2507>.

Krister Lindén and Lauri Carlson. Finnwordnet–finnish wordnet by translation. *LexicoNordica–Nordic Journal of Lexicography*, 17: 119–140, 2010.

Márton Miháلتz, Csaba Hatvani, Judit Kuti, György Szarvas, János Csirik, Gábor Prószéky, and Tamás Váradí. Methods and results of the hungarian wordnet project. In *Proceedings of The Fourth Global WordNet Conference*, pp. 311–321, 2008.

Köszönöm a figyelmet!