

Workflows for Developing Finnish–German Shallow RBMT... in FinMT, Helsinki, 2017-11-01

<https://github.com/flammie/apertium-fin-deu/>

Tommi A Pirinen (tommi.antero.pirinen@uni-hamburg.de)

HZSK.de, de.CLARIN.eu, etc.

October 27, 2017

Introduction

- ▶ A professional comp.linguist moving to Germany with “no prior knowledge” of the language
- ▶ A field/documentary linguistics-heavy North Eurasian language project (SIL FLeX!)
- ▶ Prior experience of FSAs, RBMT, SMT...
- ▶ 🖱️ RBMT+FSA+language learning+language documentation = cool new workflow for creating: corpora + mono- and bilingual lexicons, chunkers and more!!

Language documentation workflow

Field linguists / documentary linguists are used to SIL tools like FLeX (toolbox) with workflows like:

1. split words into morphs
2. tag morphs
3. translate lemmas
4. repeat

Results in a dictionary, annotated corpus and some kind of “morphological analyser”.



Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotypioilla ratsastaneen mainoksen , joka suututti jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen , mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta .



Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppioilla ratsastaneen mainoksen , joka suututti jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen , mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta .



Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkais-i meksikolaisilla stereotyyppiöilla ratsastaneen mainoksen , joka suututt-i jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkais-i viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen , mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisi-lla stereotyyppiöilla ratsastaneen mainoksen , joka suututti jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen , mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisi-lla stereotyyppiä-lla ratsastaneen mainoksen , joka suututti jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen , mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainoksen, joka suututti jopa suurlähettilään – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä.. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen, mutta vain muutaman päivän jälkeen mainos poistettiin kaikkialta.

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainoksen, joka suututti jopa suurlähettilään – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalle-en, mutta vain muutaman päivän jälkeen-en mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolais-illa stereotyyppi-illa ratsasta-ne-en mainokse-n , joka suututti jopa suurlähettilään – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalle-en , mutta vain muutaman päivän jälke-en mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokse-n , joka suututti jopa suurlähettilää-n – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viiko-n lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalle-en , mutta vain muutama-n päivä-n jälke-en mainos poistettiin kaikkialta .



Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokse-n , joka suututti jopa suurlähettilää-n – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalle-en , mutta vain muutamien päivien jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokseen, joka suututti jopa suurlähettiläään – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen, mutta vain muutamien päivien jälkeen mainos poistettiin kaikkialta.

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokseen, joka suututti jopa suurlähettilää – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalleen, mutta vain muutama päivä jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokseen, joka suututti jopa suurlähettilää-n – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatko-a suosituille Hintasaarnaaja - mainossarjalle-en, mutta vain muutama-n päivän jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokse-n , joka suututti jopa suurlähettilää-n – Nyt mainos on poistettu , eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjalle-en , mutta vain muutamien päivien jälkeen mainos poistettiin kaikkialta .



Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppiä ratsasta-ne-en mainokseen, joka suututti jopa suurlähettilää – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjoille, mutta vain muutamien päivien jälkeen mainos poistettiin kaikkialta .

Demo FLeX here (Finnish, simplified)

<https://www.hs.fi/nyt/art-2000005425359.html> Elisa julkaisi meksikolaisilla stereotyyppisillä ratsastaneiden mainosten, joka suututti jopa suurlähettilään – Nyt mainos on poistettu, eikä yhtiö halua enää edes puhua siitä .. Teleyhtiö Elisa julkaisi viime viikon lopussa jatkoa suosituille Hintasaarnaaja - mainossarjoille, mutta vain muutamien päivien jälkeen mainos poistettiin kaikkialta .
... after only 15 minutes of clicking and stuff you get like, past tense, 10 case allomorphs. passive voice, and a dozen of lemma allomorphs

Real FLeX

Kalaba - FieldWorks Language Explorer

File Edit View Data Insert Format Tools Parser Window Help

Texts & Words

- Interlinear Texts
- Concordance
- Word List Concordance
- Word Analyses
- Bulk Edit Wordforms
- Statistics

Texts

Title

Show All

My Green Mat

Text

Title

Kal

Eng My Green Mat

Info Baseline Gloss Analyze Tagging Print View Text Chart

1.1 Word

Morphemes

Lex. Entries

Lex. Gloss

Lex. Gram. Info.

Word Gloss

Word Cat.

pus

pus

pus₁

green

adj

green

mod

yalola

yalola -la

yalola -la

mat 1SgPoss

N (I) N:(Possessor)

my mat

N

nihimbilira

ni-	him-	*bili	-ra
ni-	hiN-	*bili	-ra
1SgSubj	3SgObj	to.see	Pres
V:(Subject)	V:Object	trans (1)	sta:Tense
I see			
V			

RBMT workflow (apertium style)

1. (try to) Translate a text
2. Add OOVs to source language dictionary, repeat 1
3. Add missing word-translations to bilingual dictionary, repeat 1
4. Add ungenerateable words to target language dictionary, repeat 1
5. Mangle mismatching lexical grammar
6. (Advanced RBMT not relevant to this presentation, magic)

Results in two morphological analyser/generators, a bilingual dictionary and a shallow comparative grammar!

RBMT workflow demo

Am S-Bahnhof Berliner Tor hat es gleich drei Einsätze der Bundespolizei gegeben. Donnerstagabend war dort ein stark angetrunkener Mann auf die Gleise gefallen und musste von den Beamten gerettet werden. Freitag um halb acht lief ein ebenfalls angetrunkener Mann auf den Gleisen Richtung Hauptbahnhof. Und gegen zehn warf ein Mann den Deckel eines Mülleimers auf die Gleise. In allen drei Fällen musste die Strecke **kurzzeitig** gesperrt werden.

RBMT workflow demo

Päälle *S-asema berliiniläinen portilla on se samoin kolme käytöt liittovaltionpoliisin annettu. torstai-ilta oli siellä yksi vahvasti *angetrunkener mies @Gleis<n><pl><ine> miellyttää ja täytyi virkailijoista pelastetaan. perjantai jotta puolittainen kahdeksan saapui myös *angetrunkener mies @Gleis<n><pl><nom> suunnassa @Hauptbahnhof<n><sg><nom>. Ja vastaan kymmenen pudotti mies @Deckel<n><sg><gen> @Mülleimer<n><sg><gen> @Gleis<n><pl><ine>. kaiken kolme tapauksien/tapausten täytyi reitti lyhytaikainen estetään.

RBMT workflow demo

Päälle @S-Bahnhof<n><sg><nom> berliiniläiseltä portilla on se samoin kolme käytöt liittovaltionpoliisin annettu. torstai-ilta oli siellä yksi vahvasti @angetrunken<adj><pos><sg><nom> mies @Gleis<n><pl><ine> miellyttää ja täytyi virkailijoista pelastetaan. perjantai jotta puolittainen kahdeksan saapui myös @angetrunken<adj><pos><sg><nom> mies @Gleis<n><pl><nom> suunnassa @Hauptbahnhof<n><sg><nom>. Ja vastaan kymmenen pudotti mies @Deckel<n><sg><gen> @Mülleimer<n><sg><gen> @Gleis<n><pl><ine>. kaiken kolme tapauksien/tapausten täytyi reitti lyhytaikainen estetään.

RBMT workflow demo

#lähijuna-asema<n><sg><nom> berliiniläiseltä portilla on se samoin kolme käytöt liittovaltionpoliisin annettu. torstai-ilta oli siellä yksi vahvasti humalainen mies raiteissa miellyttää ja täytyi virkailijoista pelastetaan. perjantai jotta puolittainen kahdeksan saapui myös humalainen mies raiteet suunnassa #päärautatieasema<n><sg><nom>. Ja vastaan kymmenen pudotti mies kannen roskiksen raiteissa.

RBMT workflow demo

lähijuna-asema berliiniläiseltä portilla on se samoin kolme käytöt liittovaltionpoliisin annettu. torstai-ilta oli siellä yksi vahvasti humalainen mies raiteissa miellyttää ja täytyi virkailijoista pelastetaan. perjantai jotta puolittainen kahdeksan saapui myös humalainen mies raiteet suunnassa päärautatieasema. Ja vastaan kymmenen pudotti mies kannen roskiksen raiteissa. kaiken kolme tapauksien/tapausten täytyi reitti lyhytaikainen estetään.

Experimental “results”

- ▶ 1.5 years after starting at week-daily average effort of 15 minutes
- ▶ approx. 10,000 translation pairs, few hundreds of new lexemes for analysers, few dozens of rules for WSD, chunking, lex.sel., etc.
- ▶ Post-edition WER is about 40, BLEU probably 8 or so,
- ▶ ... but it’s easy to read and understand!
- ▶ For realistic projects, e.g. within Uralic group, it should be perfectly plausible

Let’s go: <http://39476.s.time4vps.cloud>