

Fieldwork and Grammaticography in a Digital World



```
! =====
!! !Vowel rules
! =====

!! ! metaphony

"Default VH" !regressive vowel height assimilation in V1 with i/u in V2"
!! __@RULENAME@_
  Vx:Vy <=> [#].#. Cns:*= _ Cns:+ VHtrig Cns:*= Dummy:*= Vow:*= %>:0 ;
    where Vx in ( á a ä å )
          Vy in ( ä i e u )
          matched ;

"Default VH for 4syllables" !(ignores first foot, otherwise same as above"
!! __@RULENAME@_
  Vx:Vy <=> [#].#. Cns:*= Vow:+ Cns:+ Vow:+ Cns:+ _ Cns:+ VHtrig Cns:*=
    where Vx in ( á a ä å )
          Vy in ( ä i e u )
          matched ;
```

Joshua Wilbur

Freiburg Research Group in Saami Studies • Universität Freiburg

Descriptive Grammars and Typology • University of Helsinki

28 March 2019

Fieldwork and Grammaticography in a Digital World

Overview

- background
- fieldwork
- grammatographicy
- other advances
- outlook





BACKGROUND (aka: contextualization)

Pite Saami

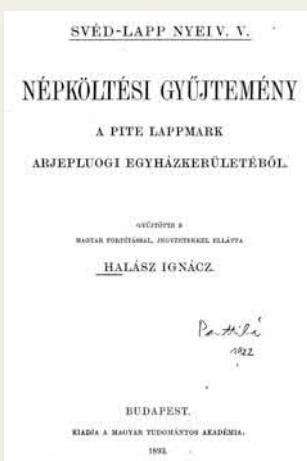
- Uralic > Finno-Ugric > Saamic
- spoken by ~40 individuals from Arjeplog/Árjepluovve in Swedish Lapland
- aka: Arjeplog-Saami, *bidumsámeigiella*
- nearly all speakers are at least 50
- all speakers are bilingual (Pite Saami and Swedish/*Arjeplogsmål*)
- no *official* orthography (yet...), but a working standard
- no media
- Swedish dominates everyday life
- hardly being passed on to younger generations



Pite Saami

larger linguistic studies:

- Halász 1893 (in Hungarian)
- Lagercrantz 1926 (in German)
- Ruong 1943 (in German)
- Lehtiranta 1992 (in Finnish)
- Wilbur 2014 (in English)
- Sjaggo 2015 (in Swedish)



Pite Saami

larger linguistic studies:

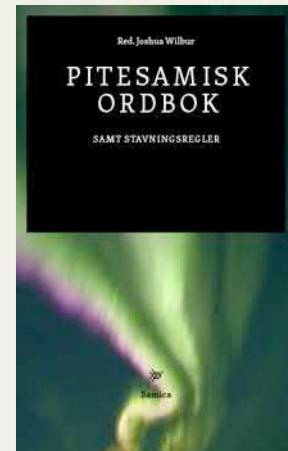
- Halász 1893 (in Hungarian)
- Lagercrantz 1926 (in German)
- Ruong 1943 (in German)
- Lehtiranta 1992 (in Finnish)
- Wilbur 2014 (in English)
- Sjaggo 2015 (in Swedish)

The screenshot shows a search interface for the Pite Saami lexical database. The search term 'Bäckgårde' is entered in the search bar. The results page displays several entries, each with a title, definition, and examples. The interface includes a sidebar with language selection (Swedish), a navigation menu, and a footer with copyright information.

other materials:

- extensive collection of heritage materials (ISOF, Uppsala)
- dictionary (Pite Saami -> Swedish/English)
and proposed orthographic rules (2016)
- online lexical database
- online orthographic rules (including spellchecker (in beta))
- smart phone app in the works

The screenshot shows a document titled 'Pitesamiska stavningsregler'. It includes sections on 'Klusiler (b, d, g, k, p, t)' and 'Konkurrens mellan klusiler och prenasillat'. The document provides examples of how specific consonants are pronounced in different contexts, such as between vowels or at the end of words. It also discusses the treatment of 'v' and 'h' in Pite Saami.



Pite Saami

larger linguistic studies:

- Halász 1893 (in Hungarian)
- Lagercrantz 1926 (in German)
- Ruong 1943 (in German)
- Lehtiranta 1992 (in Finnish)
- Wilbur 2014 (in English)
- Sjaggo 2015 (in Swedish)

other materials:

- Extensive collection of heritage materials (ISOF, Uppsala)
- dictionary (Pite Saami -> Swedish/English)
and proposed orthographic rules (2016)
- online lexical database
- online orthographic rules (including spellchecker (in beta))
- smart phone app in the works

recent linguistics projects:

- Documentation (2008-2015; materials archived at ELAR and TLA)
- Lexicography (2016)
- Syntactic structures (2016-present)

Pite Saami

larger linguistic studies:

- Halász 1893 (in Hungarian)
- Lagercrantz 1926 (in German)
- Ruong 1943 (in German)
- Lehtiranta 1992 (in Finnish)
- Wilbur 2014 (in English)
- Sjaggo 2015 (in Swedish)

other materials:

- Extensive collection of heritage materials (ISOF, Uppsala)
- dictionary (Pite Saami -> Swedish/English)
and proposed orthographic rules (2016)
- online lexical database
- online orthographic rules (including spellchecker (in beta))
- smart phone app in the works

recent linguistics projects:

- Documentation (2008-2015; materials archived at ELAR and TLA)
- Lexicography (2016)
- Syntactic structures (2016-present)

Pite Saami

-> *each fieldwork situation is unique!*

- significant aspects of mine include:
 - an accessible modern technological infrastructure on-site
 - a previous history of linguistics work
 - extensive language technology tools for closely-related languages
 - messy but extant orthographic “tradition” when I started





FIELDWORK in a digital world

tools for fieldwork

- in the old days: notebook and pencil
- nowadays:
 - recording equipment
 - laptop
 - digital backup capacity (even in the cloud)
 - transcription software (ELAN)
 - mobile phones
 - social media
 - (e.g.: for staying in contact, data source)
 - grammaticography software
 - (e.g. FLEx for interlinearization)
 - language technology
- ...

data collection and fieldwork

- modern, affordable digital recording technologies (especially video) allow fieldworkers to capture much more than just language, but the entire human event
 - more complete documentation, potentially useful beyond linguistics*

why not use:

- body cameras
- drones
- surround-sound microphones
- 360° cameras
- 3-D cameras

...

*cf. Rießler & Wilbur 2017

(re-)collecting old data (heritage harvesting)

- OCR (optical character recognition)

34

SAMEFOLKETS EGEN TIDNING

<p>Pånya kädde</p> <p>Män kalkay mujhtahit sáme, tñutuq luotah. muy mán mäddé paleu kuhiv sehtjat ahtjiet ja tñutjajiet, ku mán hivv mannam. Tat empiisse ja tñutj Vuona-lantajet, kumse Tjeldjaka sunsle falen juhtia ja kissij urrun. Ju tan tñut men tav muhtajteq tñjelat, ku tñl kirsse luvvaq pähtemin. Lij sámes sapme, mij pravje kuj kiesee jähjet Stuornon Vuona riijkan. Tak lin pär kihkuna almatja, tat same-pundi, ku ahks, ja sij siika aktav nejtau. Te lij tñ kiesee, ja kähde lij tat tan Stuornoen räkken. Ja nello kooto tñj krasse-pukkojan ja tan pakkjen jähki kattes. Te vuolkeepi, tak kuoktes vulus Vuonan vadtset, ja tat njáta pahtaa iktuk kähthaj ieloq vuoidnet. Tak kuoktes ärojka vuonan mäddé peyveh, ja ku te vist pajan vadttsika, ja ku putikja kate vujdnusini, te vuojdnepe sáj, aht iello lij kuujt kätte kuoran livan. Sáj vättisjka. Ja ku te kätte lahka putikja, mak sáj te äska åttjojka vuoid- net! Tat iello, mak sáj jähkika lij livvatem — tat lij kalle urrenn livva iello — mán mij lij sjaddam tajna ieloj ja nejtajn! Sáj vuujnike, aht njáta lij tjukhim ielov ja mannam kätte kuorraj ja livutahamt ielov ja alkana altojt pähhtjet. Ja ku lij pähjemin, te lij atja-rajte jala aija-rajte tsapmostem tav livva ielov ja pähhtjet nejtar jamas. Äpä siita iello ja almatj vil aj fieraj tasa tan kiädj.</p>	<p>iello-tjárråkav jähkä rasta. Ja te kalkaj sán ietj vist rasta mannat, men te vältij jähkä su, ja sakkjoi njáta tan jähkäj. Skajte-kietjen lij jähkä tjsaskam rumpahav kaddaj. Ja pallin almatjah räkketan skajtaj ja javestin rup- mahav tasa. Män mäddé jakeh tat rajest vässö almatjah kullin, ku tat njáta juoqkaj tan skajte-kietjen ja hälkj: »Jähkä skävvä kul mu pär pieljehtuhija.» Ja lerrin almatjah tav vuolev ja leriahttu puolvast puolvaj.</p>
<p>Suptsas atja-rajte pirra</p> <p>Tav läu mu tñjat-ajja mujhtalam. Lij sámes sapme mij iktuk urutij ja kultij pivti. Käta- tjav inij kánné vieso. Sámes palen ku lij iktuk káten, te pähnta sisä kähjat amas stuorra älma- mán nubbe juolke lä änep. Nuppen kietan lä stuorra ruovte-karre ja tidno ja nuppen kietan stuorra vähtjer. Ja te hälla sapmaj: »Ihta tal tän man tuápta?» Sapme pallaj ja vastetij aht idtijj sán taptä. Tä hälla tat stuorra älma: »Män läv tät tat atja-rajte. Ku kullapihit män läv mannamin, ta tsapmap tajna vähtjerat tan ruovte-karraj, ja ku tñdun tsapmap te altakastav, ja rassjo pähta taste, ku karest lejkhiv tñtjatse.» Ja te vit hälkj: »Ku kulla- pihit män läv pähtemin, te ehpít åttja aktakav tahkat jal parkat, ajanat verthipehtit vuorter tasak män läv mäddel mannam. Män iv vuokahä, aht tijah kalkapihitit nakanav tahkat, ku män läv mannulakan. Tav kalka tän hällit kajhka ietja almatjija.</p> <p>Ja te mán aj mujhtav, aht tulutij vuorats almatjah lin nävht mårähkhan manaja, jus idtjiu sjavot åro, ku atja-rajte jala »Balv-ajjas lij mannamin.</p>	<p>Renkalvarna fängades med metspön</p> <p>Skogssamerna från Västra Kikkejaur samlas varje sommar i Tjatjais vid Bakttive för märkning av årets renkalvar. Därvid tillämpas man ett av Gällivaresemernas prävöt system vid infangandet av kalvarna. I stället för lasso har man använt sig av långa spinn, i vars yttersta ända en repigla anbringats, och på detta sätt fångas renkalvarna i baksben.</p> <p>Två skogssamrer från Västra Kikkejaur vore för några år sedan på studieresa till Gällivare, där skogssamernas sommarrennskötsel studeras. Där ha samerna nämligen gjit in för den intensiva renskötelsen och bl. a. infört det nya ovannämnda systemet ifråga om kalvarnas infangande. Detta har bl. a. den förtjänsten, att kalvarna inte skiljas från vajorna, enär renarna förlåt sig lugnare och där kalvarna följa vajorna helt stillsamt. Framför allt har lappfoggen varit intresserad av att det nya systemet prövas, enär det har sina givna förtjänster. Gällivaresemerna ha även gjit in att följa renarna natt och dag, vilket ger renarna lugnande och trygga.</p> <p>Det var givetvis ingen litet sak för samerna att på en gång utverk från lassokastningen till smetandet, och det visar att märkningen i Tjatjais nyttjades man det gamla beprövade sättet. Renrigdet var nämligen i största laget, vilket gjorde, att smetandet blev allt för tidskrävande.</p> <p>I Tjatjais häller det på att vixa upp en lappstad i miniatyr. Där har nämligen byggts sex kåtor i stil med de, som finns i lappstaden i Arvidsjaur, två bärber och tre rengärden.</p>

Hur drängen Lars blev en stor renpatron

Följande shistoria har mina förfäder berättat för mig, och de gjorde detta så, att berättelsen fick samolikhetens prägel.

Händelsen förelägges till cirka tre kilometer norr om Araskone, där renrigdet Riebel-bävre-gidde ligger. Lars var dräng åt en rik skogssam vid namn Olof — i dagligt tal kallad Wuolla. En natt, då Lars var ute och valtade renhjorden, mötte han två näpna sameflökor — wieternäldak. De var båda klidda i sevaljaha, yttierplagg, sydda av beredda renskin, som färgats röda med den naturfärg, som framstälts av albark. Flökernas härfätor var långa och prydda med vackra pärflöband av runda pärlor. Dessutom hade dessa törser i sina långa flötor färdigtagna sensören, som samar förr i tider använde som syträd vid olika sinnader. De hade också var sin hand med sig. Den ena av flökerna frågade Lars, vilkendera han helst önskade sig: en fågelhund eller en renhund. Lars ville naturligtvis ha renhunden, som han också genast fick som gåva.

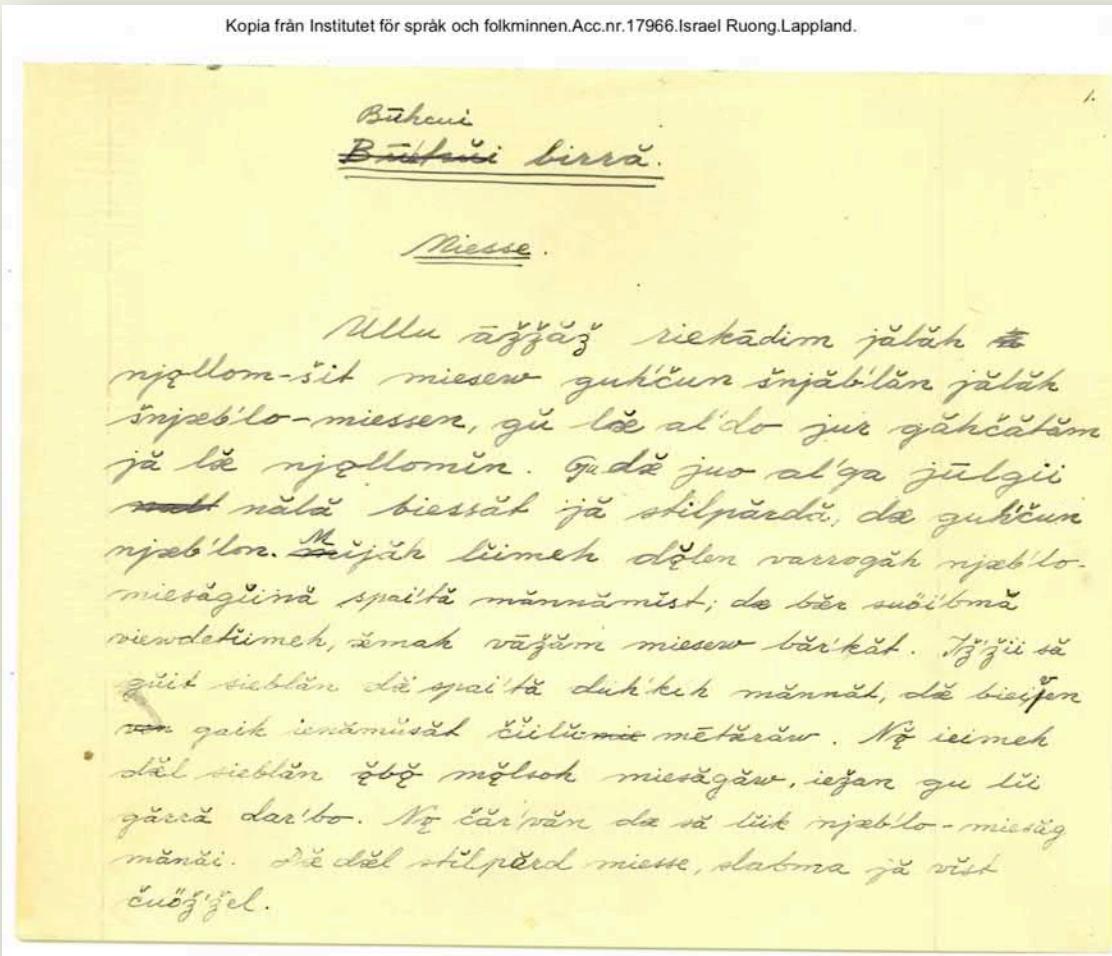
Då berättas det, att hans renar började förska sig snabbt, att han på några år blev en stor renhusbond. Till att börja med hade han visserligen ett fåtal renar, som han fätt som dränglöns under krens lopp, men sedan han blev ägare till den underbara renhunden, visste han att han blivit

- embedded text
(more than just scanning!)
- can be exported
(e.g. to ELAN)
- can be part of a corpus*

*cf. Partanen & Rießler 2019

(re-)collecting old data (heritage harvesting)

- HTR (handwritten text recognition)



- embedded text
(more than just scanning!)
- can be exported
(e.g. to ELAN)
- can be part of a corpus*
- *much more complex than OCR, thus it currently requires much more training data before it's useful*

*cf. [Transkribus](#) project (Kahle 2017);
also Blokland et al 2019 for a brief discussion

```

352 ! =====
353 !! !Vowel rules
354 !
355
356 !! ! metaphony
357
358 "Default VH" !regressive vowel height assimilation in V1 with i/u in V2" ! á:ä a:i ä:e å:u ^0:0
359 !! __@RULENAME@_
360   Vx:Vy <=> [#].#. Cns:*= _ Cns:+ VHtrig Cns:*= Dummy:* Vow:*= %>:0 ;
361   where Vx in ( á a ä å )
362     Vy in ( ä i e u )
363     matched ;
364
365 "Default VH for 4syllables" !(ignores first foot, otherwise same as above)
366 !! __@RULENAME@_
367   Vx:Vy <=> [#].#. Cns:*= Vow:+ Cns:+ Vow:+ Cns:+ _ Cns:+ VHtrig Cns:*= Dummy:* Vow:*= %>:0 ;
368   where Vx in ( á a ä å )
369     Vy in ( ä i e u )
370     matched ;

```

GRAMMATICOGRAPHY in a digital world

brief history of grammaticography

- 1/3 of the Boasian trilogy

...

- Payne 1997, Mosel 2006, Aikhenvald 2015, etc.
- Nordhoff 2008 *Electronic Reference Grammars for Typology: Challenges and Solutions*
- **Implemented grammars (incorporation in corpus and computational linguistics)**

digital tools for grammaticography

- Toolbox, FLEx
 - good for concatenative morphology
 - *play, play-s, play-ed, play-er, play-er-s*
 - not so good for non-linear morphology
 - *sing, sing-s, sang, sung*

*What do you do when
non-linear morphology is
the **default** in your
language?*

digital tools for grammaticography

- Toolbox, FLEX

What do you do when
non-linear morphology is
the **default** in your
language?

	SG	PL
NOM	<i>juällge</i>	<i>juolge</i>
GEN	<i>juolge</i>	<i>julgij</i>
ACC	<i>juolgev</i>	<i>julgijt</i>
ILL	<i>juallgáj</i>	<i>julgijda</i>
INESS	<i>juolgen</i>	<i>julgijn</i>
ELAT	<i>juolgest</i>	<i>julgijst</i>
COM	<i>julgijna</i>	<i>julgij</i>
ABESS	<i>juolgedak</i>	<i>juolgedaga</i>
ESS		<i>juallgen</i>

juällge ‘foot/leg’

digital tools for grammaticography

- Toolbox, FLEx
- other, digital approaches...

	SG	PL
NOM	<i>juällge</i>	<i>juolge</i>
GEN	<i>juolge</i>	<i>julgij</i>
ACC	<i>juolgev</i>	<i>julgijt</i>
ILL	<i>juallgáj</i>	<i>julgijda</i>
INESS	<i>juolgen</i>	<i>julgijn</i>
ELAT	<i>juolgest</i>	<i>julgijst</i>
COM	<i>julgijna</i>	<i>julgij</i>
ABESS	<i>juolgedak</i>	<i>juolgedaga</i>
ESS		<i>juallgen</i>

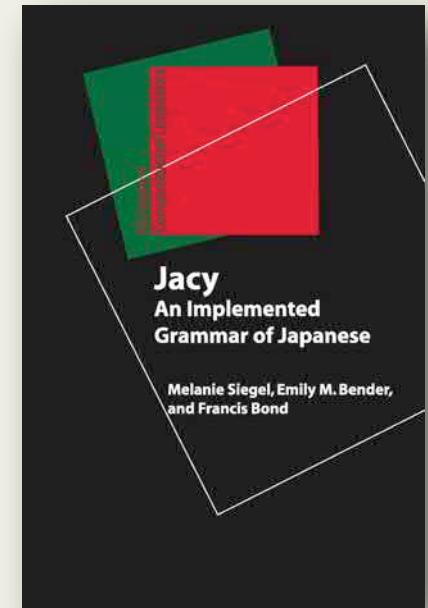
What do you do when
non-linear morphology is
the **default** in your
language?

4 stem allomorphs:
juällg-
juolg-
juallg-
julg-

juällge ‘foot/leg’

implemented grammars

- aka “precise” grammars
 - self-validating
- computer-processable
 - but only borderline human-readable
(at least from a traditionalist perspective)
 - computational linguists, typically HPSG
- **analyze linguistic structures**
- **implementation --> parse and tag a corpus**



Siegel et al. 2016

cf. new *Language Science Press* series “[Implemented Grammars](#)”

implemented grammar (FST/CG) for Pite Saami

- **Giellatekno infrastructure:**
the Research group for Saami language technology at University Tromsø
 - FST – Finite State Transducer¹
 - CG – Constraint Grammar²
- **automatic annotations in ELAN...**

¹Beesley & Karttunen 2003; ²Didriksen 2007–2018, Karlsson 1990; Karlsson et al. 1995

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

Constraint Grammar (CG) → for removing ambiguities in FST output

formalism:

- `lexc` (lexicon, PoS, linear morphology)
- `twolc` (non-linear morphology)
- `cg3` (syntax)

Uses orthographic standard!

Båñka gäddé

¶Måñ galgav mujjitalit såmes dulutj giehtov, mav måñ mådde bálen gulliv iehtjan áhtjest ja tjidtjájjást, gu måñ lijjiv mánán. Dat subtsas lä dajst Vuona ländajst, gunne tjidtjáhka såme dålen juhtin ja gesijd urrun. Ja dan dihta måñ dav mujjtájiv tjállet, gu del giesse lä ruvvaj båhtemin. Lij såmes sábme mij pruvvkuj giesen jáhtet Stuormon Vuona ríjkan. Da lin ber gålbmå almatja, dat såme bundi, suv áhkka, ja såj inijga aktav niejdav. Dä lij del giesse, ja gåhte lij del dan Stuormo vággen. Ja iello guodoj dajn grásse bahkojn ja dan vággen jåhkågáddev. Dä vuällgeba da guoktes vulus Vuonan vátset, ja dat näjja báhtsá iktuk gáhttáj ielov vuäjdnet. Da guoktes árojga Vuonan mådde biejve, ja gu dä vist bajás váttsájga, ja gu budijga gáde vujdusij, dä vuäjdneba såj, ahte iello lij gujd gáde guoran livan. Såj vättsijga. Ja gu dä gáde lahka budijga, mav såj dä ässká ådtjojga vuäjdnet? Dat iello, mav såj jähkjiga lij livvademin – dat lij galle urrum livva iello. Men mij lij sjaddam dajna ielojn ja niejdajn? Såj vujnijga, ahte näjja lij tjuhkkim ielov ja mannam gáde guorraj ja livudahttám ielov ja állgám áldojd båhtjet. Ja gu lij båhtjemin, dä lij áttjárájjde jala ájjárájjde tsábmestam dav livva ielov ja båhtjejniejdav jámas. Åbbå sijda iello ja almatj vil aj fieraj dasa dan gäddáj.

¶Ja ie lam báhtsám ienap buhtsu viessot, gu da ma lin ullgulun guohtomin jala «skånårdemin» – nåv gukkte mijá várén pruvvkujin hållåt, ja da ie lam ienap gu nagan nälljelåk hägga.

¶Dan guakktásij sjaddag lüssis vájes ja umårredis giessebiejve. Näjja lij håhkkånam ja gákja iello láttkanam, ahte såj iebá máhttám jáhtet. Ja dat rájest guhttjun dav gieddev Båñka Gádden, man nala lij åbbå sijda iello ja näjja áttjárájdest tsábmestuvvum jámas.

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

formalism:

- `lexc` (lexicon, PoS, linear morphology)
- `twolc` (non-linear morphology)

Output analyses:

```
juällge
juällge juällge+N+Sg+Nom

juallgáj
juallgáj          juällge+N+Sg+Ill

julgijd
julgijd juällge+N+Pl+Acc

juolgen
juolgen juällge+N+Sg+Ine
```

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

input:

wordform

output:

wordform lemma+PoS+Morphology

juällge
juällge

juällge+N+Sg+Nom

julgijd
julgijd

juällge+N+Pl+Acc

juällge juällge+N+Sg+Nom

juallgáj
juallgáj

juällge+N+Sg+Ill

julgijd
julgijd

juällge+N+Pl+Acc

juolgen
juolgen

juällge+N+Sg+Ine

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

formalism:

- lexc (lexicon, PoS, linear morphology)

```
611 jupptsa:jupptsa N EVEN "gruel, porridge, soup"
612 juällge:ju~llge N EVEN "leg, foot" ; ! no. 94
613 juällgebielle:juällge#bielle N EVEN "one leg (as opposed to both legs)" ; ! no. 4919
614 juällgeblárre:juällge#blárre N EVEN "the ball of the foot" ; ! no. 944
615 juällgetjuvvde:juällge#tjuvvde N EVEN "toe" ; ! no. 947
616 jáde+N+Der/NomAg:jádedieddje N EVEN "leader, chairperson" ; ! no. 90556
617 jágnjå:jágnjå N EVEN "lingonberry" ; ! no. 950
618 jähkå:jähkå N EVEN "stream, creek" ; ! no. 3434
619 já...1...jä...1... N EVEN "Christmas" ; ! no. 963

35 LEXICON N_EVEN ! giella, guolle, bissti
36 +N: EVENCASE ; ! Sg Nom, Ess
37 +N: N_EVEN_ILL ;
38 +N: ^WG N_EVEN_WK ; ! CG here
39 +N: N_EVEN_DIM ; !DIM derivation

40 LEXICON N_EVEN_J
41 +Pl+Gen: K ;
42 +Sg+Gen: K ;
43 N_ILL: K ;
44 +CG: K ;
45 +CG_WK: K ;
46 +CG_HERE: K ;
47 +CG_DERIVATION: K ;
48 +CG_NOM: K ;
49 +CG_ESS: K ;
50 +CG_ILL: K ;
51 +CG_WK: K ;
52 +CG_HERE: K ;
53 +CG_DERIVATION: K ;
54 +CG_NOM: K ;
55 +CG_ESS: K ;
56 LEXICON N_EVEN_J
57 +Pl+Gen: K ;
58 +Pl+Acc:d K ;
59 +Pl+Ine:n K ;
60 +Pl+Ela:st K ; ! norm
61 +Pl+Ela+Ilse/NG:s K ; ! norm
```

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

formalism:

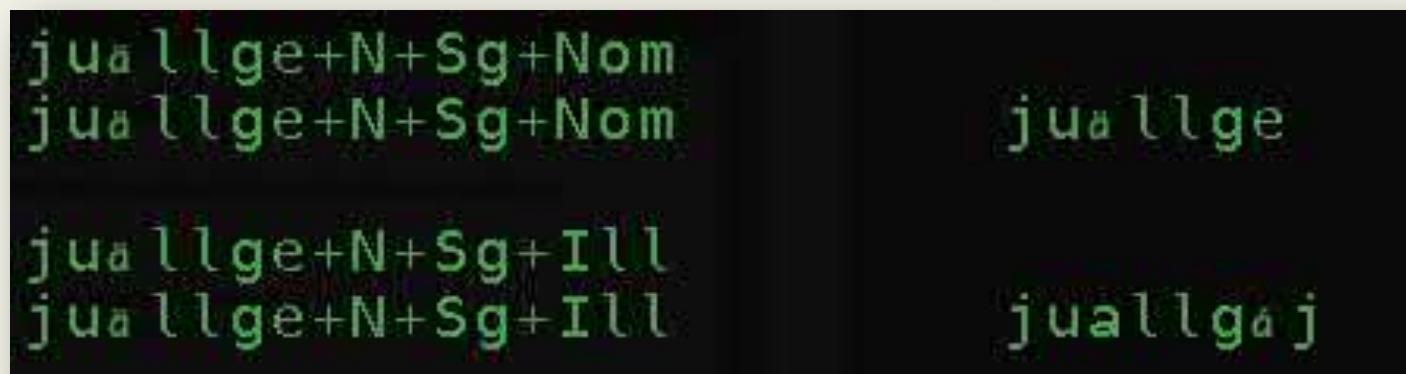
- twolc (non-linear morphology)

```
3 "Consonant Gradation for xxy:xy" ←
4 !! __@RULENAME__
5 Cx:0 <=> Vow:+ Cx Cy Vow:+ Cns:*( Dummy: ) %^WG: ;
6 where Cx in ( l l l l l l l l m m m m )
7     Cy in ( j b d f g k p s v b g p s )
8     matched ;
9
10 "Default VH" ! regressive vowel height assimilation in V1 with i/u in V2" ! á:ä a:i ä:e å:u ^0:0
11 !! __@RULENAME__
12 Vx:Vy <=> [#].#. Cns:*_ - Cns:+ VHtrig Cns:*( Dummy:*) Vow:*>:0 ;
13 where Vx in ( á a ä å )
14     Vy in ( ä i e u )
15     matched ;
16
17 "Default UA in G3" ! u^0:ua (juallgáj, luskkta); always with word-initial C
18 !! __@RULENAME__
19 %^0:a <=> Cns:*u - G3 [:a|:o|:y|:ä|:å]:á \%^WG: ;
20     Cns:*u - noCG [:a|:o|:y|:ä|:å]:á Cns:*_ %^G3: *_ %^UAUML: ;
21
22 "Special UÄ (VH) in G3" ! u^0:ua (juällge); always with word-initial C
23 !! __@RULENAME__
24 %^0:ä <=> Cns:*u - G3 :e \%^WG: ;
25     Cns:*u - Cns:+ :e (j:j) %^G3: ; ! duädde<->duodde, guäddej<->guäddeja
26     Cns:*u - noCG :e Cns: *_ %^UAUML: ; ! tjuädtjelit (not tjuodtjelit)
27
28 "Special VH for u^0" ! u^0:u0
29 !! __@RULENAME__
30 %^0:0 <=> Cns:*u - Cns:+ VHtrig Cns:*( Dummy:*) %>:0 ; ! juällge<->julgij
31
32 "V2 E to I before j-suffixes" ! guolle -> gulij
```

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for **generating** wordforms



(it works in both directions)

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

```
julgij  
julgij  juallge+N+Pl+Gen  
julgij  juallge+N+Pl+Com  
  
julgijn  
julgijn juallge+N+Pl+Ine  
julgijn juallge+N+Sg+Com  
  
gieße  
gieße  giesse+N+Pl+Nom  
gieße  giesse+N+Sg+Gen  
gieße  giesset+V+ConNeg  
gieße  giesset+V+Imprt+Sg2  
gieße  giesset+V+Vgen
```

BUT: how to deal with **morphologically ambiguous** wordforms?
(disambiguation)

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

Constraint Grammar (CG) → for removing ambiguities in FST output

formalism:

- lexc (lexicon, PoS, linear morphology)
- twolc (non-linear morphology)
- cg3 (syntax)

```
328 #NP (incl. pronouns) preceding certain Po is Gen
329 SELECT:genB4Po Gen IF (*1C Po BARRIER NoNP);
330 REMOVE:NoPoWithoutGen Po IF (NOT *-1 Gen BARRIER NoNP);
331 #NP (incl. pronouns) following certain Pr is Gen
332 SELECT:genAfterPr Gen IF (*-1C Pr BARRIER NoNP);
333 REMOVE:NoPrWithoutGen Pr IF (NOT *1 Gen BARRIER NoNP);
```

example: rules describing dependency between adpositions and genitive case

implemented grammar (FST/CG) for Pite Saami

infrastructure:

Finite State Transducer (FST) → for analyzing wordforms

Constraint Grammar (CG) → for removing ambiguities in FST output

formalism:

- lexc (lexicon, PoS, linear morphology)
- twolc (non-linear morphology) **output** (analyses)
- cg3 (syntax)

```
328 "m&an"
      "m&an" Pron Pers Sg1 Nom
329 "<lev>" 
      "l&a" V Ind Prs Sg1
330 "<julgij>" 
      "ju&llge" N Pl Gen SELECT:329:genB4Po
331 "#NP (incl. pro)
      ;   "ju&llge" N Pl Com SELECT:329:genB4Po
332 SELECT:genAfterNP
      ;   "nanne" Po
333 REMOVE:NoPrWithNP
      ;   "nanne" Adv REMOVE:313:NoAdvIfPo
```

disambiguation example

nala

gähettjat

tjurvij

daj

disambiguation example

nala

gähttjat

tjurvij

daj

onto

look+INF

antler+GEN+PL

DET+GEN+PL

antler+COM+PL

DET+COM+PL

PRON+GEN+PL

PRON+COM+PL

FST output:

nala

nala nala+Po

gähttjat

gähttjat gähttjat+V+Inf

tjurvij

tjurvij tjårrve+N+Pl+Gen

tjurvij tjårrve+N+Pl+Com

daj

daj dat+Det+Pl+Gen

daj dat+Det+Pl+Com

daj dat+Pron+Dem+Pl+Gen

daj dat+Pron+Dem+Pl+Com

disambiguation example

daj

tjurvij

nala

gähettjat

‘to look at those antlers’

[pit100405b.011]

FST output:

```
nala
nala      nala+Po

gähettjat
gähettjat      gähettjat+V+Inf

tjurvij
tjurvij  tjårrve+N+Pl+Gen
tjurvij  tjårrve+N+Pl+Com

daj
daj      dat+Det+Pl+Gen
daj      dat+Det+Pl+Com
daj      dat+Pron+Dem+Pl+Gen
daj      dat+Pron+Dem+Pl+Com
```

disambiguation example

<i>daj</i>	<i>tjurvij</i>	<i>nala</i>	<i>gähettjat</i>
da-j	tjurvi-j	nala	gähettja-t
DET-GEN.PL	antler-GEN.PL	onto	look-INF
'to look at those antlers'			[pit100405b.011]

FST output:

```
nala
nala      nala+Po

gähettjat
gähettjat      gähettjat+V+Inf

tjurvij
tjurvij  tjårrve+N+Pl+Gen
tjurvij  tjårrve+N+Pl+Com

daj
daj      dat+Det+Pl+Gen
daj      dat+Det+Pl+Com
daj      dat+Pron+Dem+Pl+Gen
daj      dat+Pron+Dem+Pl+Com
```

CG syntactic disambiguation:

- postpositions govern genitive NPs
`SELECT Gen IF (*1C Po BARRIER NoNP);`
- pronouns are not embedded in an NP
`REMOVE Pron IF (*1C N BARRIER NPNH);`

implemented grammars

pros:

- entirely digital (easy copying, versioning, etc.)
- computer-processable
- can analyze AND generate (useful for practical tools, e.g. teaching apps)
- accuracy can be tested on *real empirical* data
- prose can be included (as `<!--comments-->`)
- **further use in other, digital applications...**

cons:

- requires significant technical knowhow to learn and to implement
- not very human-readable, especially for non-specialists
 - prose is only included as `<!--comments-->`
 - not ideal for standard average typologists
 - not even close to ideal for most non-linguists

further use in other, digital applications...

- spell-checkers
- grammar-checkers
- teaching materials (e.g. apps)

Johan Lasko ja suv áhka Hingga viessomájge birra

The life history of Johan Lasko and his wife Hingga

...

sje19210000a-lagercrantz1957a-426

by Maria Johansson, originally transcribed by Eliel Lagercrantz in 1921.

The Finno-Ugric transcriptions and German translations were originally published as text 426 in Eliel Lagercrantz, 1957. *West- und südlappische Texte. Lappische Volksdichtung, volume 1. Suomalais-ugrilaisen Seuran Toimituksia 112. Helsinki: Suomalais-Ugrilainen Seura*; these are published by permission of Suomalais-Ugrilainen Seura.

All other transcriptions and analyses copyright Joshua Wilbur and the Pite Saami Documentation Project, licensed by CC BY-NC-SA 4.0. Lexical, wordclass and morphological analyses as well as English glosses were derived automatically, thanks in part to the Giellatekno infrastructure and to scripting assistance by Iris Perkmann and Ciprian Gerstenberger.

Johan Lasko ja suv áhka Hingga viessomájge birra.

FUT: Johan Lasko ja sù áhkà Hinka viessyöm-äjkjè pirræ.

EN: The life history of Johan Lasko and his wife Hingga.

DE: Über die Lebensgeschichte des Johan Lasko und seiner Frau Hingga.



Hingga lij riegádam Árjepluove nuortabiele suoknon, jáben 1844, ja suv ålmáj aj.

FUT: Hingga lij ri · èkäram árja-pluvjè nöö · nñta-pi:èljè su:oknuon, jápjén 1844, ja sù olmaj áj.

EN: Hinga was born on the north side of Arjeplog parish in 1844, and her husband as well.

DE: Hinga war im J. 1844 im Nordteil der Gemeinde Arjeplog geboren und so auch ihr Mann.

Sán lij nuorap ietjas áhkast

FUT: soñlij nù · òrap i · èçgs áhkäst.

EN: He was younger than his wife.

DE: Dieser war jünger als seine Frau.

áhkka

N Sg Ela

grandmother, wife

further use in other, digital applications...

- spell-checkers
 - grammar-checkers
 - teaching materials (e.g. apps)
- ...
- in documentary linguistics / endangered language descriptions
 - automatic tokenization and annotation for corpora

further use in other, digital applications...

- tier structure in ELAN corpora (Freiburg-style)

The screenshot shows the ELAN 5.0 interface with the following components:

- Video Preview:** A small video frame showing a person standing in a snowy forest.
- Annotation List:** A table showing annotations from 22.022 to 30.030. The selected row (27.028) is highlighted in purple and contains the text "tjähppis båtsoj ja".
- Timeline:** A horizontal timeline at the bottom showing time points from 00:01:52.000 to 00:01:53.400.
- Tier Structure Tree:** A hierarchical tree view on the right side, color-coded by tier type:
 - ref@AEF [96]:** Root node.
 - orth@AEF [96]:** Child of ref@AEF, containing "tjähppis båtsoj ja".
 - ft-eng@A [94]:** Child of orth@AEF, containing "black reindeer and".
 - ft-swe@A [94]:** Child of orth@AEF, containing "svart ren och".
 - word@A [473]:** Child of ref@AEF, containing "tjähppis" and "båtsoj".
 - lemma [493]:** Child of word@A, containing "tjähppat" and "båtsoj".
 - pos@[501]:** Child of word@A, containing "A" and "N".
 - morp [642]:** Child of pos@[501], containing "Attr" and "reindeer".
 - gloss [501]:** Child of pos@[501], containing "black".

Annotations for the selected tier:

- ref@AEF [96]:** 0.028 (highlighted in yellow)
- orth@AEF [96]:** tjähppis båtsoj ja
- ft-eng@A [94]:** black reindeer and
- ft-swe@A [94]:** svart ren och
- word@A [473]:**
 - tjähppis
 - båtsoj
- lemma [493]:**
 - tjähppat
 - båtsoj
- pos@[501]:**
 - A
 - N
- morp [642]:**
 - Attr
 - reindeer
- gloss [501]:** black

including annotations for:

- Lemma
- Part of speech
- Morphological categories
- Gloss

further use in other, digital applications...

- tier structure in ELAN corpora (Freiburg-style)

corpus building/extension using a **script**¹ that:

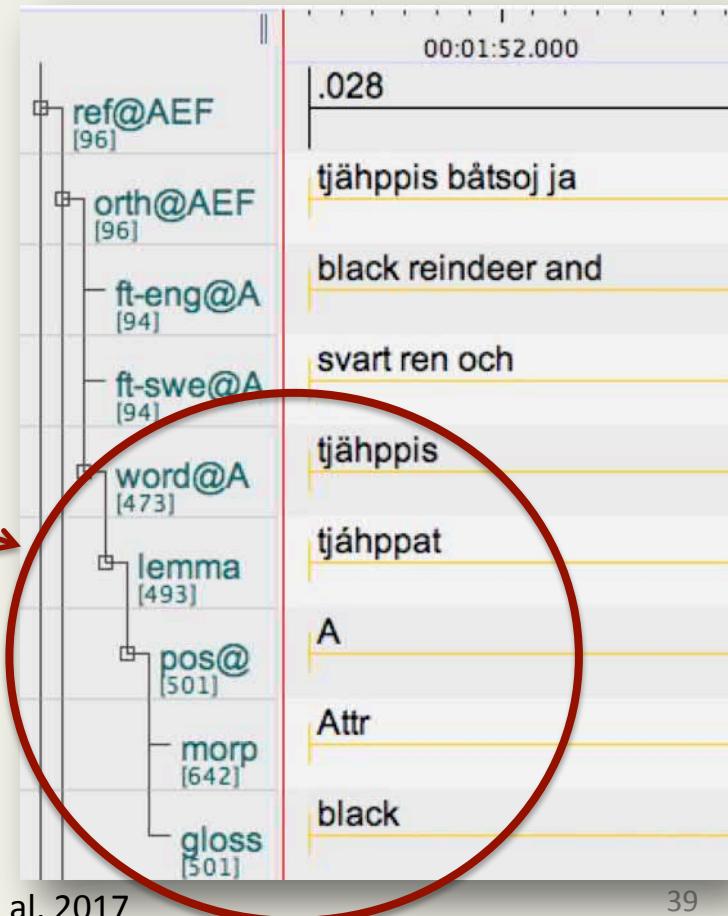
1. tokenizes the orthographic representation
2. sends each token through FST
3. removes ambiguities using CG
4. adds an English gloss
5. inserts this output into ELAN

benefits:

- saves time
- avoids inconsistencies
- can be updated automatically

*More details in talk at 11:30 in room 13
by Blokland, Partanen and Rießler*

¹cf. Blokland et al. 2015; Gerstenberger et al. 2016; Gerstenberger et al. 2017



summary of digital grammaticography

requires:

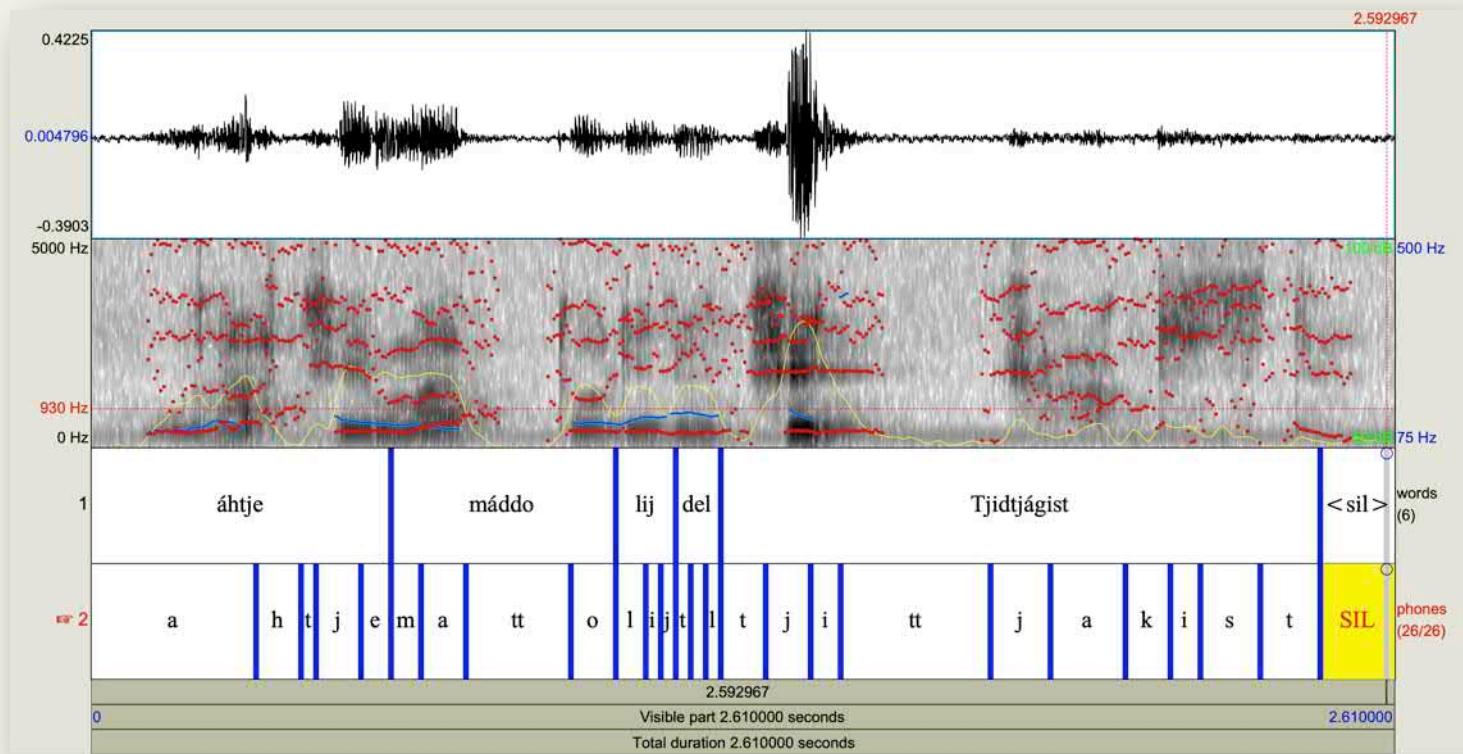
- time to learn the formalism and set up the infrastructure
- understanding of grammatical structures
- string-based representation of language

main benefits:

- can be freely accessible online
- possibility to publish (hopefully getting academic recognition, cf. LangSciPress series)
- export data for use in other tools and disciplines
 - spell-checker
 - lexicographic materials (including smart phone apps)
 - corpus building
 - teaching materials
 - increased status for the language
 - more accessible to other disciplines, e.g. via text search

main drawbacks:

- not terribly human-accessible
- not taught traditionally in General Linguistics programs



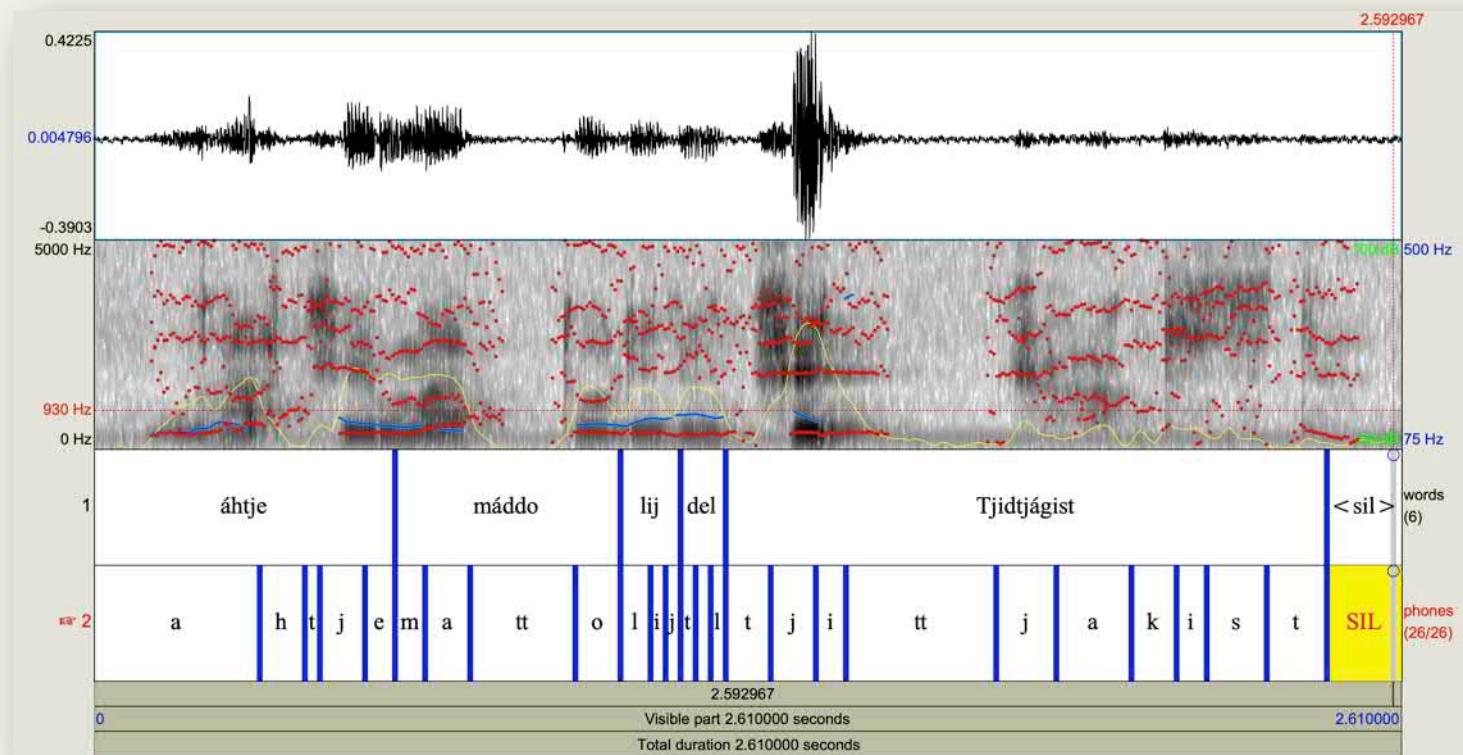
OTHER ADVANCES in digital technologies

new language technologies

- automatic segmentation, e.g.:

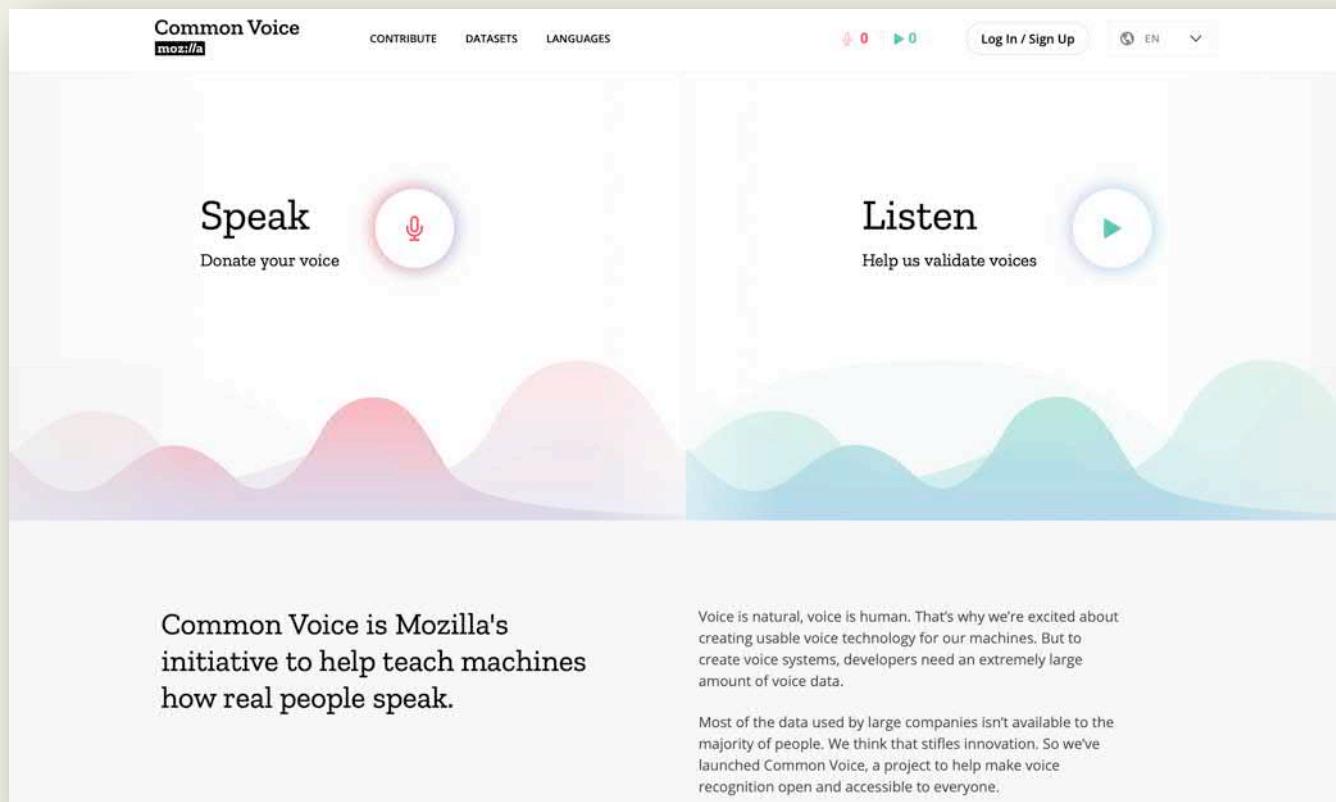
- Autosegmenteerija 2.0

- Estonian autosegmentation forced-alignment tested on Pite Saami with surprisingly accurate results:



new language technologies

- speech recognition, e.g.:
 - [Common Voice](#) (moz://a) in community development for a number of smaller languages (e.g.: Erzya, Komi-Zyrian, ...)



new language technologies

- automatic implemented grammar production
 - LinGO Grammar Matrix

<http://matrix.ling.washington.edu/customize/matrix.cgi>

LinGO Grammar Matrix

Matrix customization and download page [documentation]

Version of Tue Nov 13 21:13:22 UTC 2018

The [LinGO Grammar Matrix](#) is developed at the University of Washington in the context of the [DELPH-IN Consortium](#), by [Emily M. Bender](#) and co-authors. This material is based up work supported by the National Science Foundation under Grant No. BCS-0644097. Additional support for Grammar Matrix development has been provided by a gift to the Turing Center from the Utilika Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. The Grammar Matrix customization system is hosted by the University of Washington.

[\[University of Washington Website Terms and Conditions of Use\]](#)

[\[University of Washington Online Privacy Statement\]](#)

Publications reporting on work based on grammars derived from this system should cite [Bender, Flickinger and Oepen 2002 \[bib\]](#) and [Bender et al. 2006 \[bib\]](#). Further publications from the project are available on the [project website](#).

Filling out this form will produce a starter grammar for a natural language, consisting of a language-independent core and customized support for the language of interest. The grammar can be used as a starting point for further development. The system is experimental and may contain errors. Be advised that this system is highly experimental. We are interested in your feedback. If you have questions or comments, please email the author(s) at: ebender@u.washington.edu.

[\[Back to Matrix main page\]](#)

NOTE: Throughout the questionnaire, questions or subpages that lack a required answer or contain an incorrect answer are marked with a red asterisk (*). Subpages that contain answers that might be problematic, but are not outright incorrect, are marked with a red question mark (?). Hovering the mouse over a red asterisk or question mark will show a tooltip describing the error. Clicking on a red asterisk or question mark that is on the main page will link to the corresponding subpage.

- ▶ * [General Information](#)
- ▶ * [Word Order](#)
- ▶ [Number](#)
- ▶ * [Person](#)
- ▶ [Gender](#)
- ▶ * [Case](#)
- ▶ [Adnominal Possession](#)
- ▶ [Direct-inverse](#)
- ▶ [Tense, Aspect and Mood](#)

- ▶ * [General Information](#)
- ▶ * [Word Order](#)
- ▶ [Number](#)
- ▶ * [Person](#)
- ▶ [Gender](#)
- ▶ * [Case](#)
- ▶ [Adnominal Possession](#)
- ▶ [Direct-inverse](#)
- ▶ [Tense, Aspect and Mood](#)
- ▶ [Evidentials](#)
- ▶ [Other Features](#)
- ▶ [Sentential Negation](#)
- ▶ [Coordination](#)
- ▶ [Matrix Yes/No Questions](#)
- ▶ [Information Structure](#)
- ▶ [Argument Optionality](#)
- ▶ [Nominalized Clauses](#)
- ▶ [Clausal Complements](#)
- ▶ [Clausal Modifiers](#)
- ▶ ? [Lexicon](#)
- ▶ [Morphology](#)
- ▶ [Import Toolbox Lexicon](#)
- ▶ [Test Sentences](#)
- ▶ [Test by Generation Options](#)

Archive type: .tar.gz .zip

[Create Grammar](#) [Test by Generation](#)

new language technologies

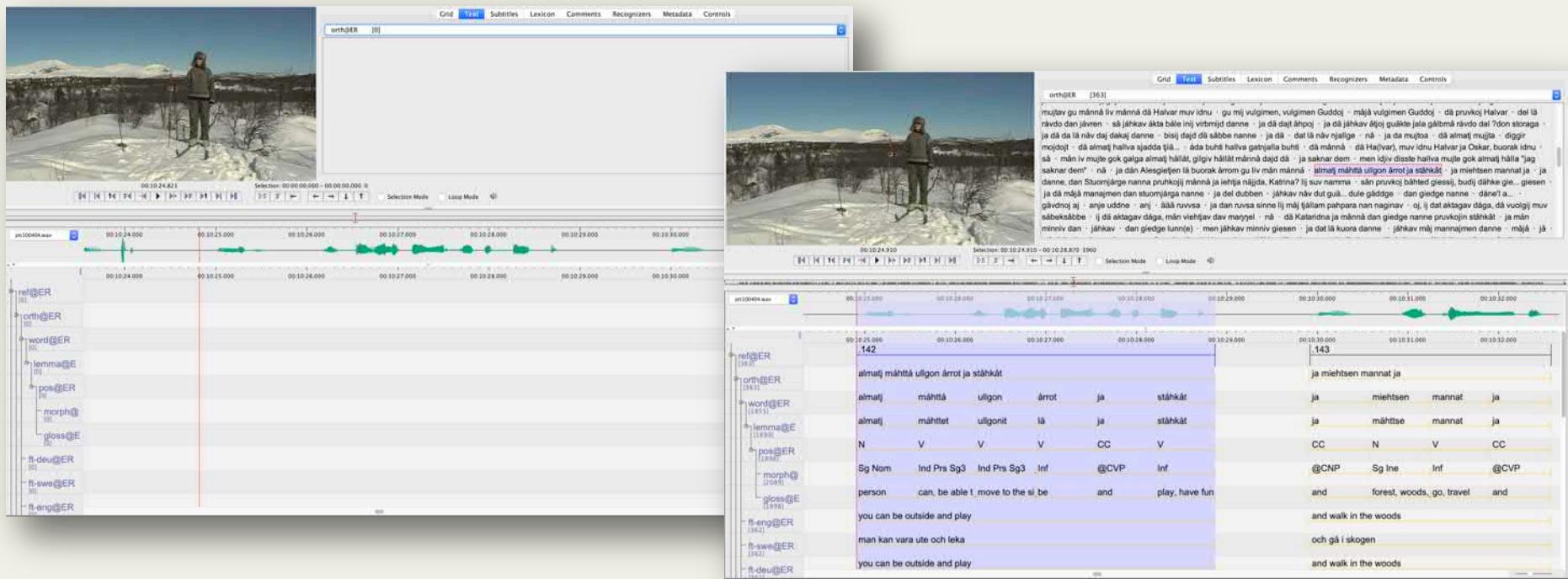
- automatic implemented grammar production
 - LinGO Grammar Matrix

<http://matrix.ling.washington.edu/customize/matrix.cgi>

```
Li 1 ;;; -*- Mode: TDL; Coding: utf-8 -*-  
M 2 ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;  
Ve 3 ;;; Grammar of mini English  
4 ;;; created at:  
5 ;;;     Wed Mar 27 12:19:05 UTC 2019  
6 ;;; based on Matrix customization system version of:  
7 ;;;     Tue Nov 13 21:13:22 UTC 2018  
8 ;;;  
9 ;;; This is a sample choices file which exercises only a small range of  
10 ;;; the information provided by the customization system, in order to  
11 ;;; create a grammar for a very small fragment of English. It describes  
12 ;;; and SVO language with a small vocabulary drawn from English and  
13 ;;; subset of the (already simple) English verbal agreement paradigm.  
14 ;;; Where it was not possible to leave a section blank, we have said  
15 ;;; the language does not manifest the phenomenon, even when this is  
16 ;;; not actually correct for English.  
17 ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;; ;;  
18 ; Type assigning empty mod list. Added to basic types for nouns, verbs and determiners.  
19  
20 non-mod-lex-item := lex-item &  
21     [ SYNSEM.LOCAL.CAT.HEAD.MOD < > ].  
22 ;;; ;;; ;;; ;;; ;;  
23 ;;; Matrix Type Addenda  
24 ;;; ;;; ;;; ;;; ;;; ;;
```

new speech technologies

- relevant technologies being developed continuously
- leading to a significant increase in efficiency for corpus building
 - > *better grammatical descriptions*





OUTLOOK

outlook

- digital tools can provide powerful advantages for both fieldwork and (especially) grammaticography and documentation
- *but*: they require knowhow that goes beyond a typical linguist's training
- I'm not saying this is for everyone, and realistically only parts will be relevant for a few – the point is:
Digital technologies should be considered, too!

References

- Aikhenvald, Alexandra Y. (2015). *The art of grammar. A practical guide*. Oxford: Oxford University Press.
- Beesley, Kenneth R. & Lauri Karttunen (2003). *Finite State Morphology*. Stanford: Center for the Study of Language and Information.
- Blokland, Rogier, Ciprian Gerstenberger, Marina Fedina, Niko Partanen, Michael Rießler, & Joshua Wilbur (2015). "Language documentation meets language technology". In: *First International Workshop on Computational Linguistics for Uralic Languages, 16th January, 2015, Tromsø, Norway. Proceedings of the workshop*. Ed. by Tommi A. Pirinen, Francis M. Tyers, & Trond Trosterud. Septentrio Conference Series 2015:2. Tromsø: The University Library of Tromsø, pp. 8–18.
- Blokland, Rogier, Niko Partanen, Michael Rießler, & Joshua Wilbur (2019). "Using computational approaches to integrate endangered language legacy data into documentation corpora. Past experiences and challenges ahead". In: *Proceedings of the Workshop on Computational Methods for Endangered Languages*. Vol. 2. Honolulu: Association for Computational Linguistics, pp. 24–30.
- Didriksen, Tino (2007–2018). *Constraint grammar manual. 3rd version of the CG formalism variant*. GrammarSoft ApS.
- Gerstenberger, Ciprian, Niko Partanen, Michael Rießler, & Joshua Wilbur (2016). "Utilizing language technology in the documentation of endangered Uralic languages". In: *Northern European Journal of Language Technology* 4, pp. 29–47.
- Gerstenberger, Ciprian, Niko Partanen, Michael Rießler, & Joshua Wilbur (2017). "Instant annotations. Applying NLP methods to the annotation of spoken language documentation corpora". In: *International Workshop on Computational Linguistics for Uralic languages (IWCLUL 2017)*. Ed. by Tommi A. Pirinen, Michael Rießler, Trond Trosterud, & Francis M. Tyers. St. Petersburg: Association for Computational Linguistics, pp. 25–36.
- Halász, Ignácz (1893). *Népköltési gyűjtemény. A Pite Lappmark arjepluogi egyházkerületéből*. Vol. 5. Svéd-Lapp Nyelv. Budapest: Magyar tudományos akadémia.
- Kahle, Philip, Sebastian Colutto, Günter Hackl, & Günger Mühlberger (2017). "Transkribus. A Service Platform for Transcription, Recognition and Retrieval of Historical Documents". In: *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 04, pp. 19–24.
- Karlsson, Fred (1990). "Constraint Grammar as a framework for parsing unrestricted text". In: *Proceedings of the 13th International Conference of Computational Linguistics*. Ed. by Hans Karlsgren. Vol. 3. Helsinki, pp. 168–173.
- Karlsson, Fred, Atro Voutilainen, Juha Heikkilä, & Arto Anttila, eds. (1995). *Constraint Grammar. A language-independent system for parsing unrestricted text*. Natural Language Processing 4. Berlin: Mouton de Gruyter.
- Lagercrantz, Eliel (1926). *Sprachlehre des Westlappischen nach der Mundart von Arjeplog*. Suomalais-ugrilaisen Seuran Toimituksia 55. Helsinki: Suomalais-Ugrilainen Seura.
- Lehtiranta, Juhani (1992). *Arjeploginsaamen öänne- ja taivutusopin pääpiirteet*. Suomalais-ugrilaisen Seuran toimituksia 212. Helsinki: Suomalais-Ugrilainen Seura.
- Mosel, Ulrike (2006). "Grammaticography. The art and craft of writing grammars". In: *Catching language. The standing challenge of grammar writing*. Ed. by Felix Ameka, Alan Dench, & Nicholas Evans. Trends in linguistics: studies and monographs 167. Berlin: Mouton de Gruyter, pp. 41–68.
- Nordhoff, Sebastian (2008). "Electronic Reference Grammars for Typology: Challenges and Solutions". In: *Language Documentation and Conservation* 2.2, pp. 296–324.
- Partanen, Niko & Michael Rießler (2019). "An OCR system for the Unified Northern Alphabet". In: *International Workshop on Computational Linguistics for Uralic languages (IWCLUL 2019)*. Tartu: Association for Computational Linguistics, pp. 77–89.
- Payne, Thomas E. (1997). *Describing morphosyntax. A guide for field linguists*. Cambridge: Cambridge University Press.
- Rießler, Michael & Joshua Wilbur (2017). "Documenting endangered oral histories of the Arctic. A proposed symbiosis for language documentation and oral history research, illustrated by Saami and Komi examples". In: *Oral history meets linguistics*. Ed. by Erich Kasten, Katja Roller, & Joshua Wilbur. Exhibitions and Symposia. Fürstenberg: Kulturstiftung Sibirien, pp. 31–64.
- Ruong, Israel (1943). *Lappische Verbalableitung dargestellt auf Grundlage des Pitelappischen*. Uppsala: Almqvist och Wiksell.
- Siegel, Melanie, Emily M. Bender, & Francis Bond (2016). *Jacy. An Implemented Grammar of Japanese*. CSLI Studies in Computational Linguistics. Stanford: CSLI Publications.
- Sjaggo, Ann-Charlotte (2015). *Pitesamisk grammatik. en jämförande studie med lulesamiska*. Senter for samiske studiers skriftserie 20. Tromsø: Septentrio Academic Publishing.
- Wilbur, Joshua (2014). *A grammar of Pite Saami*. Studies in Diversity Linguistics 5. Berlin: Language Science Press.
- Wilbur, Joshua, ed. (2016). *Pitesamisk ordbok samt stavningsregler*. Samica 2. Freiburg: Albert-Ludwigs-Universität Freiburg.



Gijtov adnet!

<i>gijtov</i>	<i>adnet</i>
<i>gijto-v</i>	<i>adne-t</i>
<i>thank-ACC.SG</i>	<i>have-PL.IMP</i>

*with special thanks to
Michael Rießler, Niko Partanen, Rogier Blokland and Ciprian Gerstenberger
for ideas, collaboration and inspiration*

Joshua Wilbur
Pite Saami Syntax Project
Freiburg Research Group in Saami Studies
joshua.wilbur@skandinavistik.uni-freiburg.de