

Strategic targeting of rich inflectional morphology for linguistic analysis and L2 acquisition / TROLLing for linguistic data

Laura A. Janda, UiT The Arctic University of Norway

CLEAR

Cognitive Linguistics: Empirical Approaches to Russian



Members of the SMARTool team



Radovan Bast



Tore Nesset



Francis Tyers



Mikhail Kopotev



Valentina Zhukova



Elizaveta Kibisova



Svetlana Sokolov



Evgeniia Sudarikova



Ekaterina Rakhilina



Olga Lyashevskaya



James McDonald

Financing and collaboration

Financed by the Norwegian Agency for International Cooperation and Quality Enhancement in Higher Education as a collaboration between UiT the Arctic University of Norway and the Higher School of Economics in Moscow

CLEAR

Cognitive Linguistics: Empirical Approaches to Russian

Diku



Overview

- **Strategic targeting of rich inflectional morphology for linguistic analysis and L2 acquisition**
 - Challenges of rich inflectional morphology
 - Zipfian distribution
 - “Paradigm Cell Filling Problem”
 - A data-based model of rich morphology
 - A machine-learning experiment
 - The SMARTool
- **TROLLing for linguistic data**
 - A professionally-curated archive for linguists everywhere

Strategic targeting of rich inflectional morphology for linguistic analysis and L2 acquisition

- **What I used to believe:**

- students need to memorize full paradigms in order to achieve fluency

- **What I now believe:**

- full paradigms are a fiction and we can strategically target acquisition of morphology

Challenges of rich inflectional morphology

- Most languages have rich inflectional morphology
- Rich inflectional morphology can involve huge numbers of forms
- Inflectional morphology can be complex

Most languages have rich inflectional morphology: A sample from Europe

HIGH	Slavic: Bosnian, Croatian, Czech, Montenegrin, Polish, Russian, Rusyn, Serbian, Slovak, Slovene, Sorbian Other Indo-European: Albanian, Greek, Irish, Latvian, Lithuanian, Romany, Scottish Gaelic, Welsh Uralic : Estonian, Finnish, Hungarian, Kven, Saami Other: Basque, Maltese, Turkish
MID	Slavic: Bulgarian, Macedonian Other Indo-European: Catalan, French, German, Italian, Luxembourgish, Portuguese, Spanish, Yiddish
LOW	Germanic språk: Danish, Dutch, English, Norwegian, Swedish

Rich inflectional morphology can involve huge numbers of forms

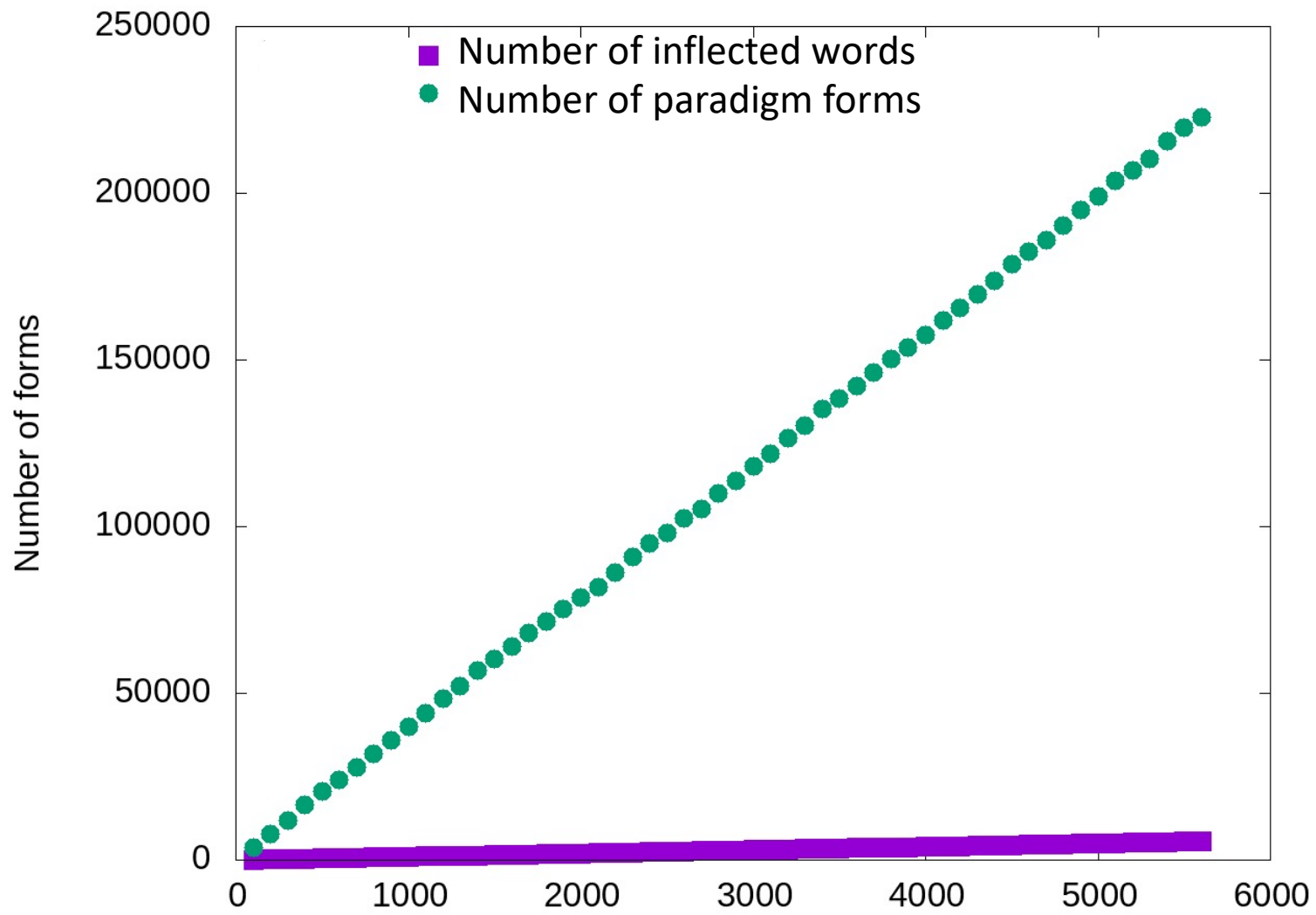
Kolik řečí znáš, tolikrát jsi člověkem

'How many languages you know – that many times you are a person'

Word	Gloss	Grammatical categories marked by morphology
<i>Kolik</i>	How many	Indefinite numeral in Accusative case
<i>řečí</i>	language	Feminine noun in Genitive case and Plural number
<i>znáš</i>	know	Imperfective verb in Present tense, Second person Singular
<i>tolikrát</i>	that many times	Adverb
<i>jsi</i>	be	Imperfective verb in Present tense, Second person Singular
<i>člověkem</i>	person	Masculine noun in Instrumental case and Singular number

Inflectional morphology in a Czech proverb

In a language with rich inflectional morphology, even a small basic vocabulary can entail **hundreds of thousands of paradigm forms.**



Inflectional morphology can be complex

- Morphology is considered both to be essential to L2 acquisition and to be a “bottleneck”, more difficult than both syntax and semantics (Slabakova 2009 & 2014, Jensen et al. 2019)
- Morphophonemic alternations can complicate the picture – inflectional morphology is not always a matter of adding desinences to stems
 - consonant and vowel alternations
 - fleeting vowels (jers)
 - suprasegmental alternations

North Saami example

Basic paradigm of <i>guoibmi</i> “partner”	
NOM.SG	<i>guoibmi</i>
GEN.SG=ACC.SG	<i>guoimmi</i>
ILL.SG	<i>guoibmá-i</i>
LOC.SG	<i>guoimmi-s</i>
COM.SG=LOC.PL	<i>guimmi-in</i>
NOM.PL	<i>guoimmi-t</i>
GEN.PL=ACC.PL	<i>guimmi-id</i>
ILL.PL	<i>guimmi-ide</i>
COM.PL	<i>guimmi-iguin</i>
ESS	<i>guoibmi-n</i>

Table 2. 81 additional paradigm forms required by NPx for *guoibmi* “partner”

NOM.SG:		GEN.SG=ACC.SG:		ILL.SG:	
1SG	<i>guoibmá-n</i>	1SG	<i>guoibmá-n</i>	1SG	<i>guoibmá-s-an</i>
2SG	<i>guoibmá-t</i>	2SG	<i>guoimmá-t</i>	2SG	<i>guoibmá-s-at</i>
3SG	<i>guoibmi-s</i>	3SG	<i>guoimmi-s</i>	3SG	<i>guoibmá-s-is</i>
1DU	<i>guoibmá-me</i>	1DU	<i>guoibmá-me</i>	1DU	<i>guoibmá-s-eame</i>
2DU	<i>guoibmá-de</i>	2DU	<i>guoimmá-de</i>	2DU	<i>guoibmá-s-eatte</i>
3DU	<i>guoibmi-ska</i>	3DU	<i>guoimmi-ska</i>	3DU	<i>guoibmá-s-easkka</i>
1PL	<i>guoibmá-met</i>	1PL	<i>guoibmá-met</i>	1PL	<i>guoibmá-s-eamet</i>
2PL	<i>guoibmá-det</i>	2PL	<i>guoimmá-det</i>	2PL	<i>guoibmá-s-eattet</i>
3PL	<i>guoibmi-set</i>	3PL	<i>guoimmi-set</i>	3PL	<i>guoibmá-s-easet</i>
LOC.SG:		COM.SG=LOC.PL:		GEN.PL=ACC.PL (=NOM.PL 1SG/DU/PL):	
1SG	<i>guoimmi-st-an</i>	1SG	<i>guimmi-in-an</i>	1SG	<i>guimmi-id-an</i>
2SG	<i>guoimmi-st-at</i>	2SG	<i>guimmi-in-at</i>	2SG	<i>guimmi-id-at</i>
3SG	<i>guoimmi-st-is</i>	3SG	<i>guimmi-in-is</i>	3SG	<i>guimmi-id-is</i>
1DU	<i>guoimmi-st-eame</i>	1DU	<i>guimmi-in-eame</i>	1DU	<i>guimmi-id-eame</i>
2DU	<i>guoimmi-st-eatte</i>	2DU	<i>guimmi-in-eatte</i>	2DU	<i>guimmi-id-eatte</i>
3DU	<i>guoimmi-st-easkka</i>	3DU	<i>guimmi-in-easkka</i>	3DU	<i>guimmi-id-easkka</i>
1PL	<i>guoimmi-st-eamet</i>	1PL	<i>guimmi-in-eamet</i>	1PL	<i>guimmi-id-eamet</i>
2PL	<i>guoimmi-st-eattet</i>	2PL	<i>guimmi-in-eattet</i>	2PL	<i>guimmi-id-eattet</i>
3PL	<i>guoimmi-st-easet</i>	3PL	<i>guimmi-in-easet</i>	3PL	<i>guimmi-id-easet</i>
ILL.PL:		COM.PL:		ESS:	
1SG	<i>guimmi-idas-an</i>	1SG	<i>guimmi-id-an-guin</i>	1SG	<i>guoibmi-n-an</i>
2SG	<i>guimmi-idas-at</i>	2SG	<i>guimmi-id-at-guin</i>	2SG	<i>guoibmi-n-at</i>
3SG	<i>guimmi-idas-as</i>	3SG	<i>guimmi-id-is-guin</i>	3SG	<i>guoibmi-n-is</i>
1DU	<i>guimmi-idas-ame</i>	1DU	<i>guimmi-id-eame-guin</i>	1DU	<i>guoibmi-n-eame</i>
2DU	<i>guimmi-idas-ade</i>	2DU	<i>guimmi-id-eatte-guin</i>	2DU	<i>guoibmi-n-eatte</i>
3DU	<i>guimmi-idas-aska</i>	3DU	<i>guimmi-id-easkka-guin</i>	3DU	<i>guoibmi-n-easkka</i>
1PL	<i>guimmi-idas-amet</i>	1PL	<i>guimmi-id-eamet-guin</i>	1PL	<i>guoibmi-n-eamet</i>
2PL	<i>guimmi-idas-adet</i>	2PL	<i>guimmi-id-eattet-guin</i>	2PL	<i>guoibmi-n-eattet</i>
3PL	<i>guimmi-idas-aset</i>	2PL	<i>guimmi-id-easet-guin</i>	3PL	<i>guoibmi-n-easet</i>

Zipfian distribution

Language & Corpus Name	Corpus Size	Paradigm Size	Total Lexemes	Lexemes with full Paradigm	% Lexemes with full Paradigm
English Web Treebank	254,830	2	6,369	1,524	23.92%
Norwegian Dependency Treebank	311,277	4	12,587	393	3.12%
Russian SynTagRus	1,032,644	12	21,945	13	0.06%
Czech Prague Dependency Treebank	1,509,242	14	17,904	3	0.02%
Estonian ArborEst	234,351	28	14,075	0	0%

Zipfian distribution

Language & Corpus Name	Corpus Size	Paradigm Size	Total Lexemes	Lexemes with full Paradigm	% Lexemes with full Paradigm
English ArborEst	254,331	20	14,073	3,340	23.92%
English ArborEst	254,331	20	14,073	439	3.12%
English ArborEst	254,331	20	14,073	13	0.06%
English ArborEst	254,331	20	14,073	3	0.02%
English ArborEst	254,331	20	14,073	0	0%

Because Zipf's Law scales up, these numbers will never change substantially, no matter how large the corpus is

Language Exposure as a Big Corpus

- A large corpus is a close approximation to the **lifetime linguistic input** for a native speaker, estimated at about 5-10 million words per year
- Zipfian distributions **remain the same** even for very large corpora, like those that approximate a speaker's exposure to their native language
- A native speaker of Russian encounters all twelve paradigm forms of **less than 0.1% of nouns** that they are exposed to in a lifetime
- The portion of adjectives and verbs attested in all paradigm forms is **virtually zero**.



What this means for words

- At the word level, each lemma has a unique “grammatical profile” -- frequency distribution of its paradigm forms

Full paradigms

	<i>'fear'</i>	<i>'soldier'</i>	<i>'department'</i>	<i>'concept'</i>	<i>'memory'</i>
sg.nom	страх	солдат	отделение	концепция	память
sg.gen	страха	солдата	отделения	концепции	памяти
sg.dat	страху	солдату	отделению	концепции	памяти
sg.acc	страх	солдата	отделение	концепцию	память
sg.ins	страхом	солдатом	отделением	концепцией	памятью
sg.loc	страхе	солдате	отделении	концепции	памяти
pl.nom	страхи	солдаты	отделения	концепции	памяти
pl.gen	страхов	солдат	отделений	концепций	памятей
p.dat	страхам	солдатам	отделениям	концепциям	памятям
pl.acc	страхи	солдат	отделения	концепции	памяти
pl.ins	страхами	солдатами	отделениями	концепциями	памятями
pl.loc	страхах	солдатах	отделениях	концепциях	памятях

Paradigm forms attested in 1M words

	<i>'fear'</i>	<i>'soldier'</i>	<i>'department'</i>	<i>'concept'</i>	<i>'memory'</i>
sg.nom	страх	солдат	отделение	концепция	память
sg.gen	страха	солдата	отделения	концепции	памяти
sg.dat	страху	солдату	отделению	концепции	памяти
sg.acc	страх	солдата	отделение	концепцию	память
sg.ins	страхом	солдатом	отделением	концепцией	памятью
sg.loc	страхе		отделении	концепции	памяти
pl.nom	страхи	солдаты	отделения		
pl.gen	страхов	солдат	отделений	концепций	
pl.dat		солдатам			
pl.acc	страхи	солдат	отделения	концепции	
pl.ins	страхами		отделениями	концепциями	
pl.loc	страхах	солдатах	отделениях		

bold >20%, plain >10%, gray 1-9%, (blank) not attested

Full paradigms

	<i>'background'</i>	<i>'champion'</i>	<i>'expanse'</i>	<i>'frame'</i>	<i>'problem'</i>
sg.nom	фон	чемпион	протяжение	рамка	трудность
sg.gen	фона	чемпиона	протяжения	рамки	трудности
sg.dat	фону	чемпиону	протяжению	рамке	трудности
sg.acc	фон	чемпиона	протяжение	рамку	трудность
sg.ins	фоном	чемпионом	протяжением	рамкой	трудностью
sg.loc	фоне	чемпионе	протяжении	рамке	трудности
pl.nom	фоны	чемпионы	протяжения	рамки	трудности
pl.gen	фонов	чемпионов	протяжений	рамок	трудностей
p.dat	фонам	чемпионам	протяжениям	рамкам	трудностям
pl.acc	фоны	чемпионов	протяжения	рамки	трудности
pl.ins	фонами	чемпионами	протяжениями	рамками	трудностями
pl.loc	фонах	чемпионах	протяжениях	рамках	трудностях

Paradigm forms attested in 1M words

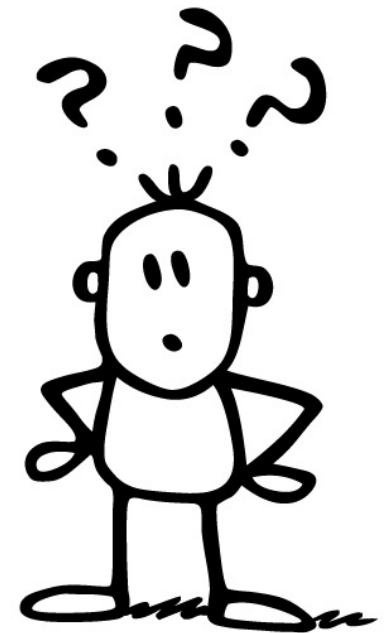
	<i>'background'</i>	<i>'champion'</i>	<i>'expanse'</i>	<i>'frame'</i>	<i>'problem'</i>
sg.nom	фон	чемпион			трудность
sg.gen	фона	чемпиона			трудности
sg.dat		чемпиону			трудности
sg.acc		чемпиона			трудность
sg.ins		чемпионом			трудностью
sg.loc	фоне		протяжении		
pl.nom		чемпионы		рамки	трудности
pl.gen		чемпионов		рамок	трудностей
p.dat		чемпионам			
pl.acc		чемпионов		рамки	трудности
pl.ins		чемпионами		рамками	трудностями
pl.loc				рамках	трудностях

bold >20%, plain >10%, gray 1-9%, (blank) not attested

Paradigm Cell Filling Problem

(Ackerman et al. 2009)

Native **speakers** of languages with **complex inflectional morphology** routinely **recognize** and **produce** forms that they have **never encountered**.



Example:

Russian gerund *nedokarmlivaja* ‘while underfeeding’ has no attestations in the Russian National Corpus (>360M words \approx lifetime exposure), but all native speakers of Russian can be expected to readily understand and to produce these forms in appropriate contexts, as evidenced by rare occurrences that turn up in Google searches.

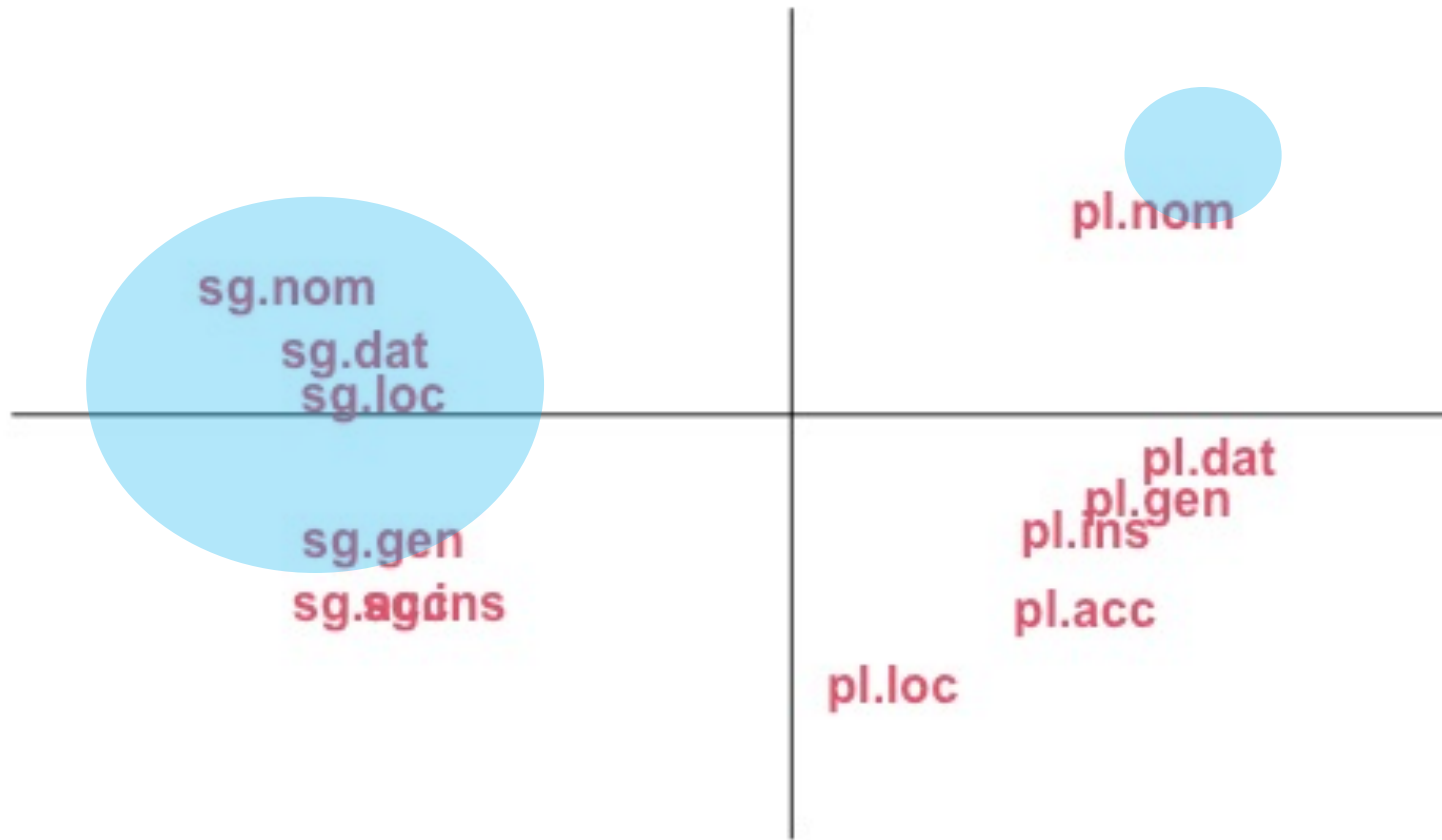
A data-based model of rich morphology

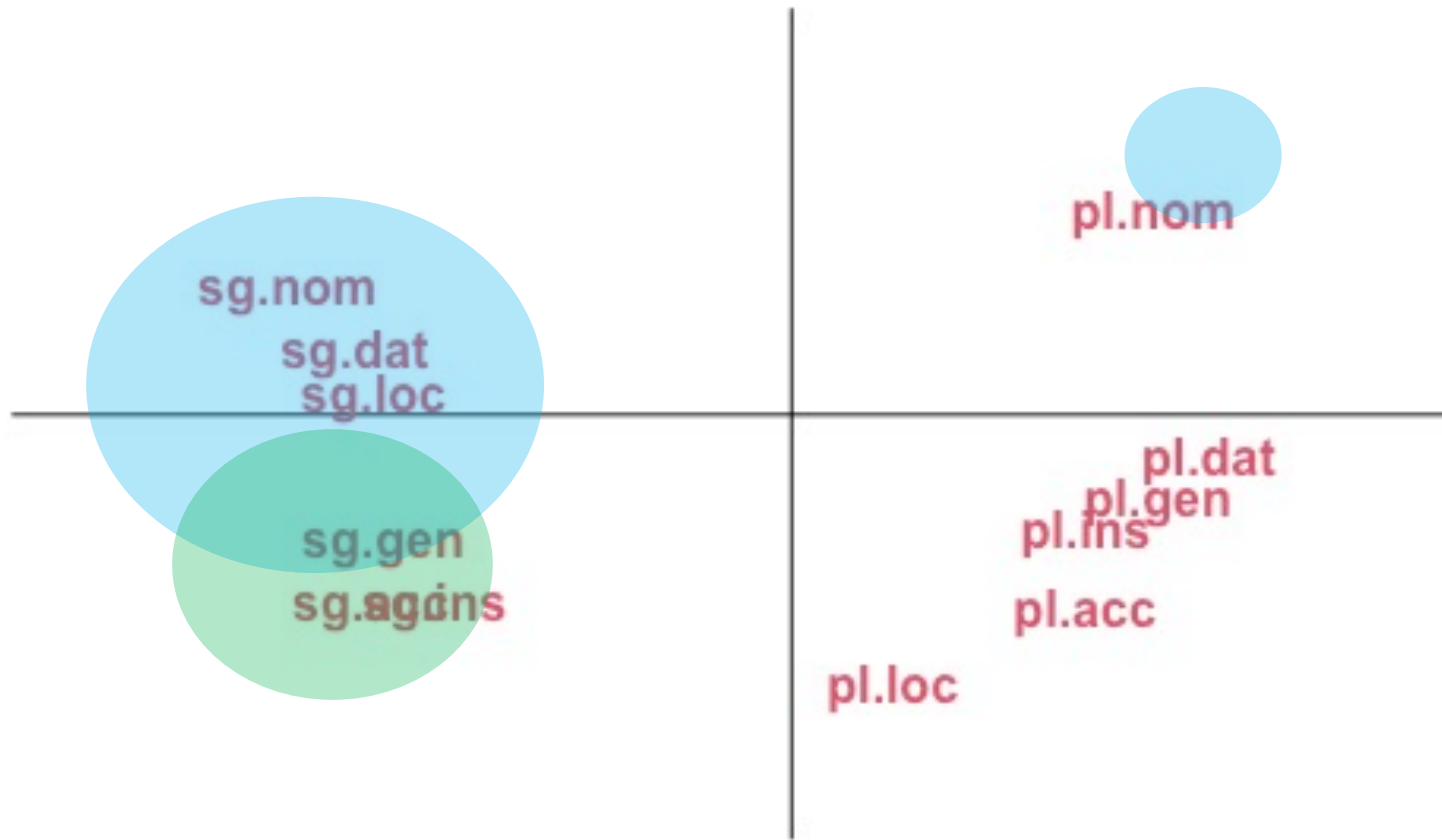
Instead of a table of equiprobable cells:

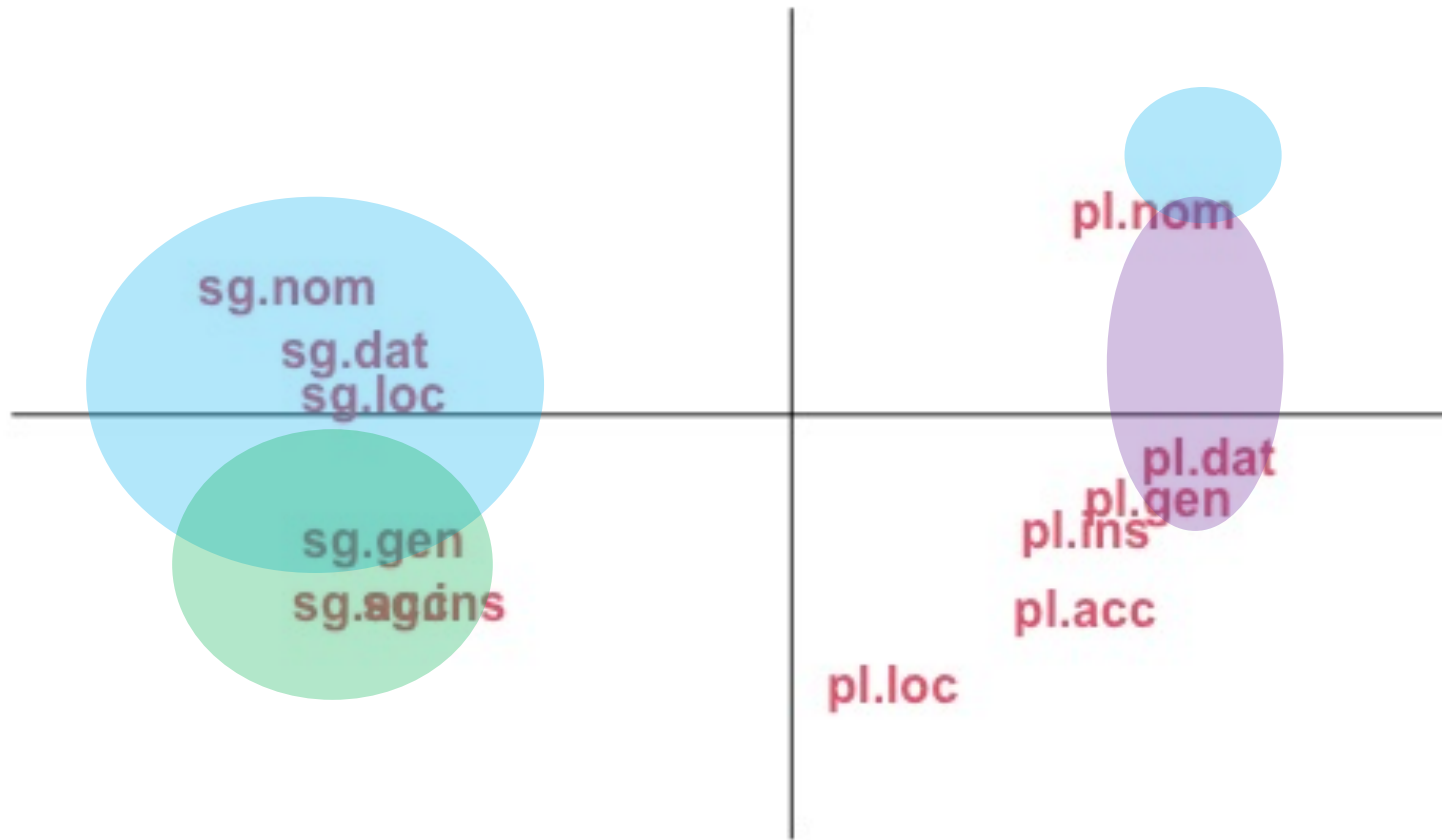
- A multidimensional space that is **partially populated in many different ways** by many different lexemes
- Inflectional morphology can be mastered through exposure to **partially overlapping subsets of paradigms**
- This **cognitively plausible** model explains why native speakers have the intuition that full paradigms exist

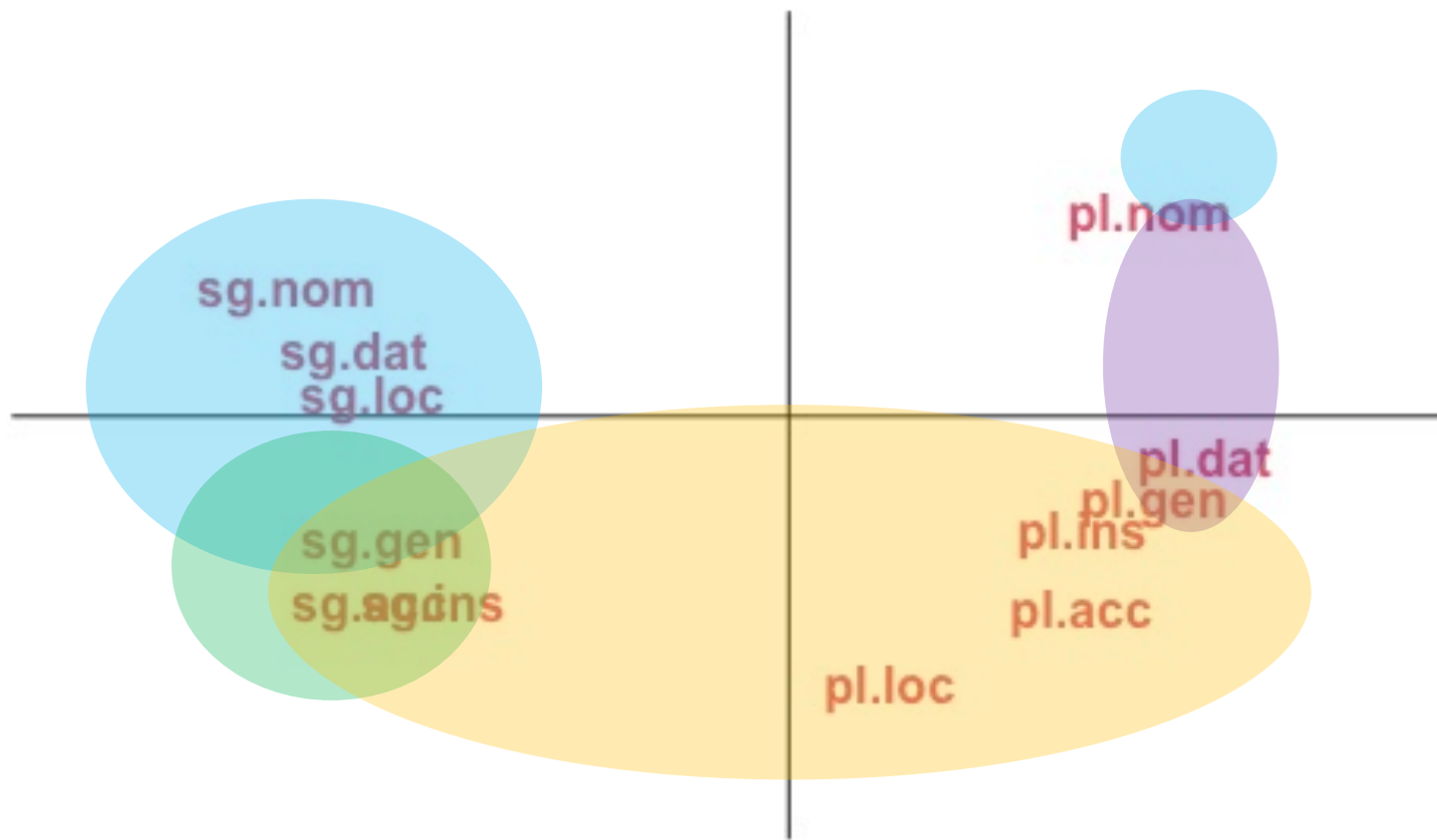
The “space” of the paradigm of Russian masculine animate nouns as determined by corpus data

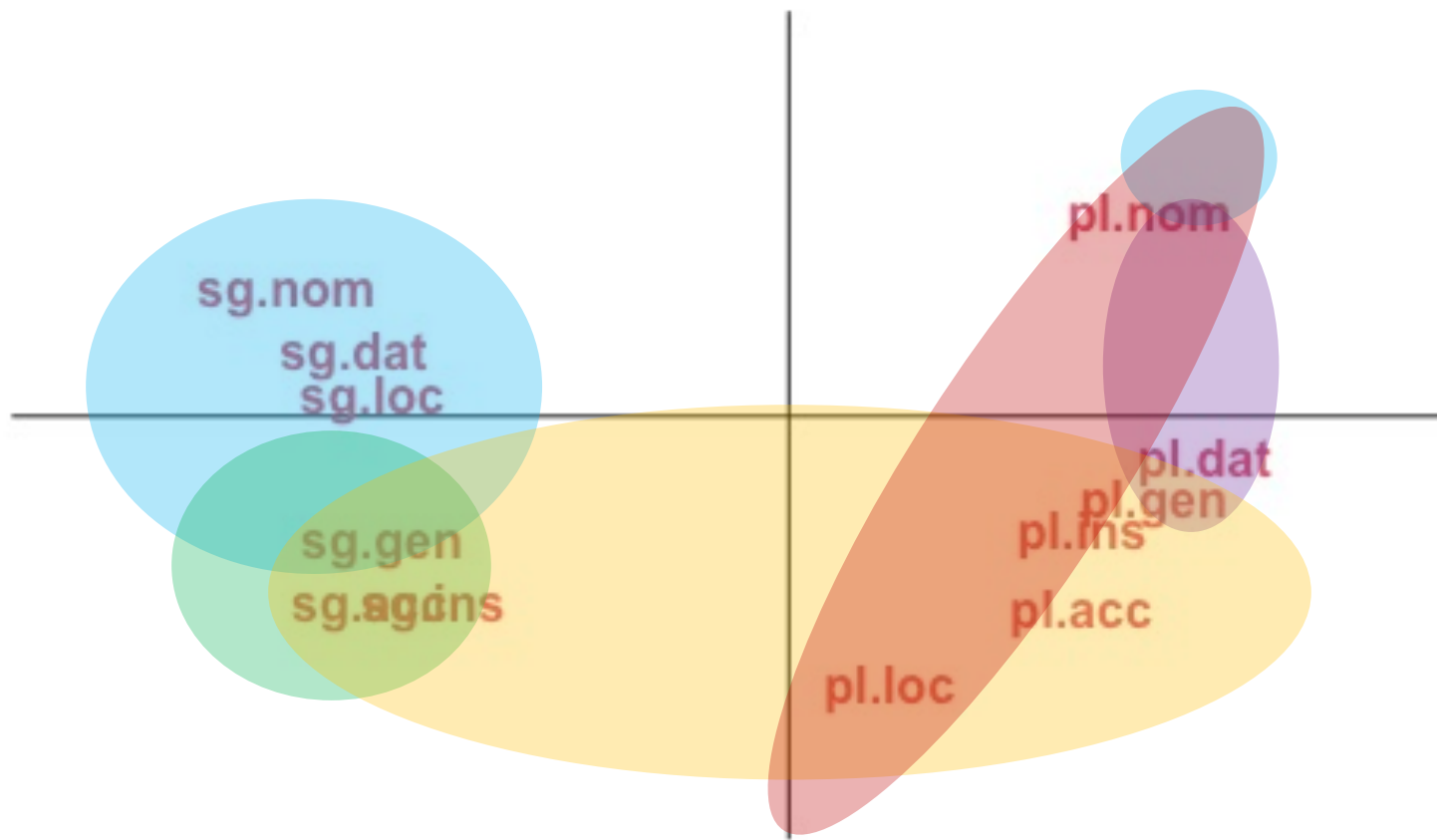


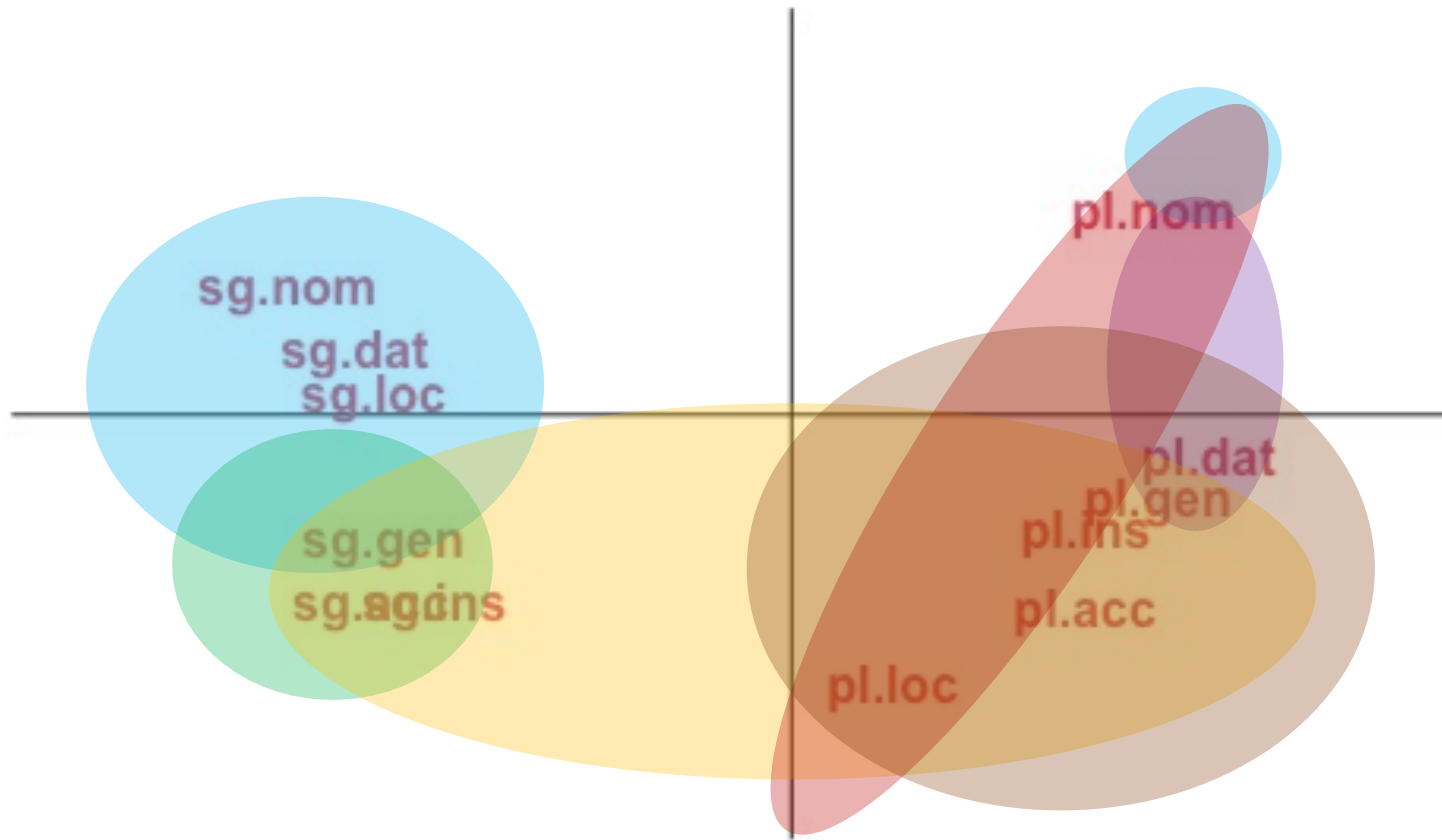


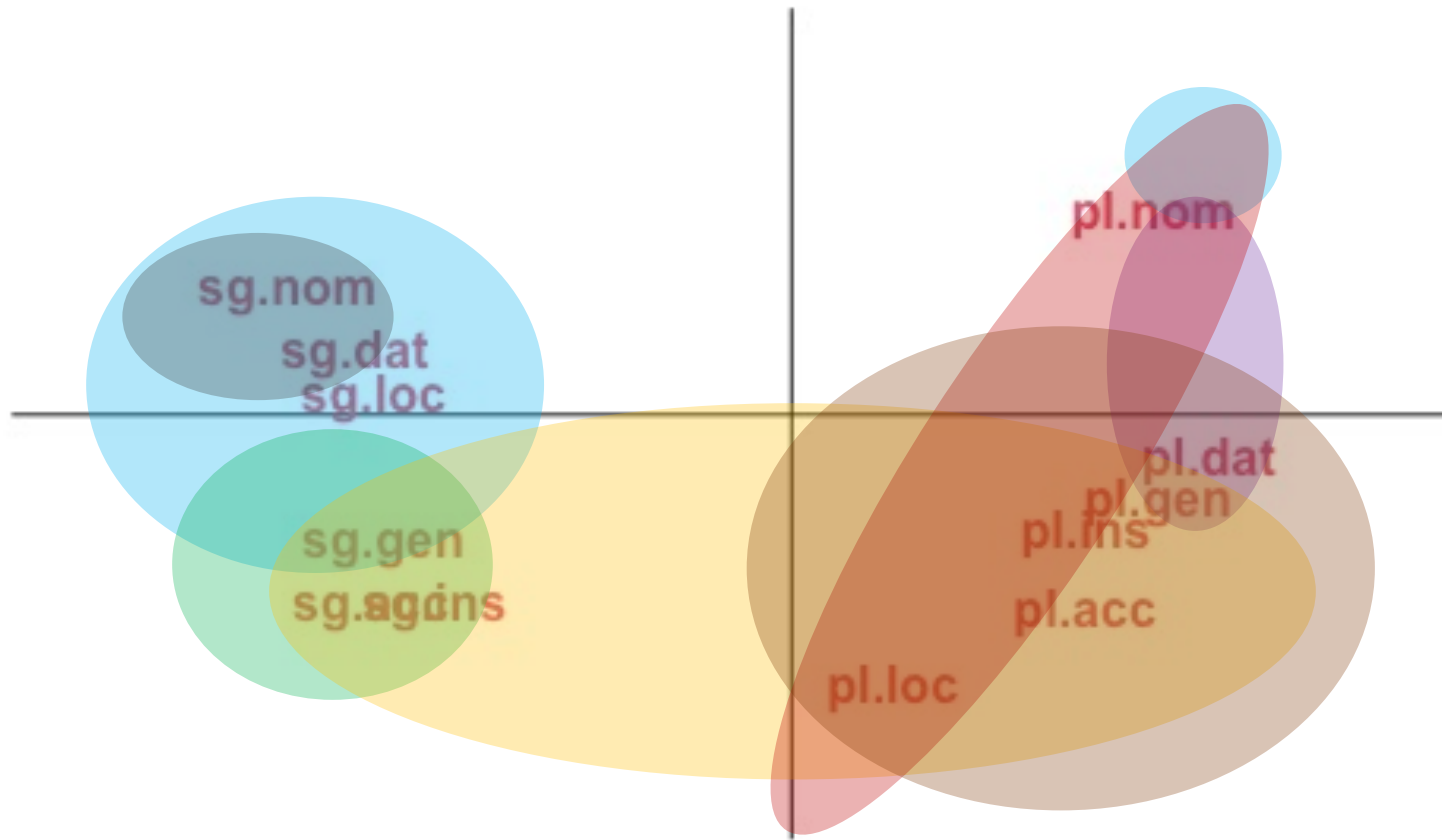


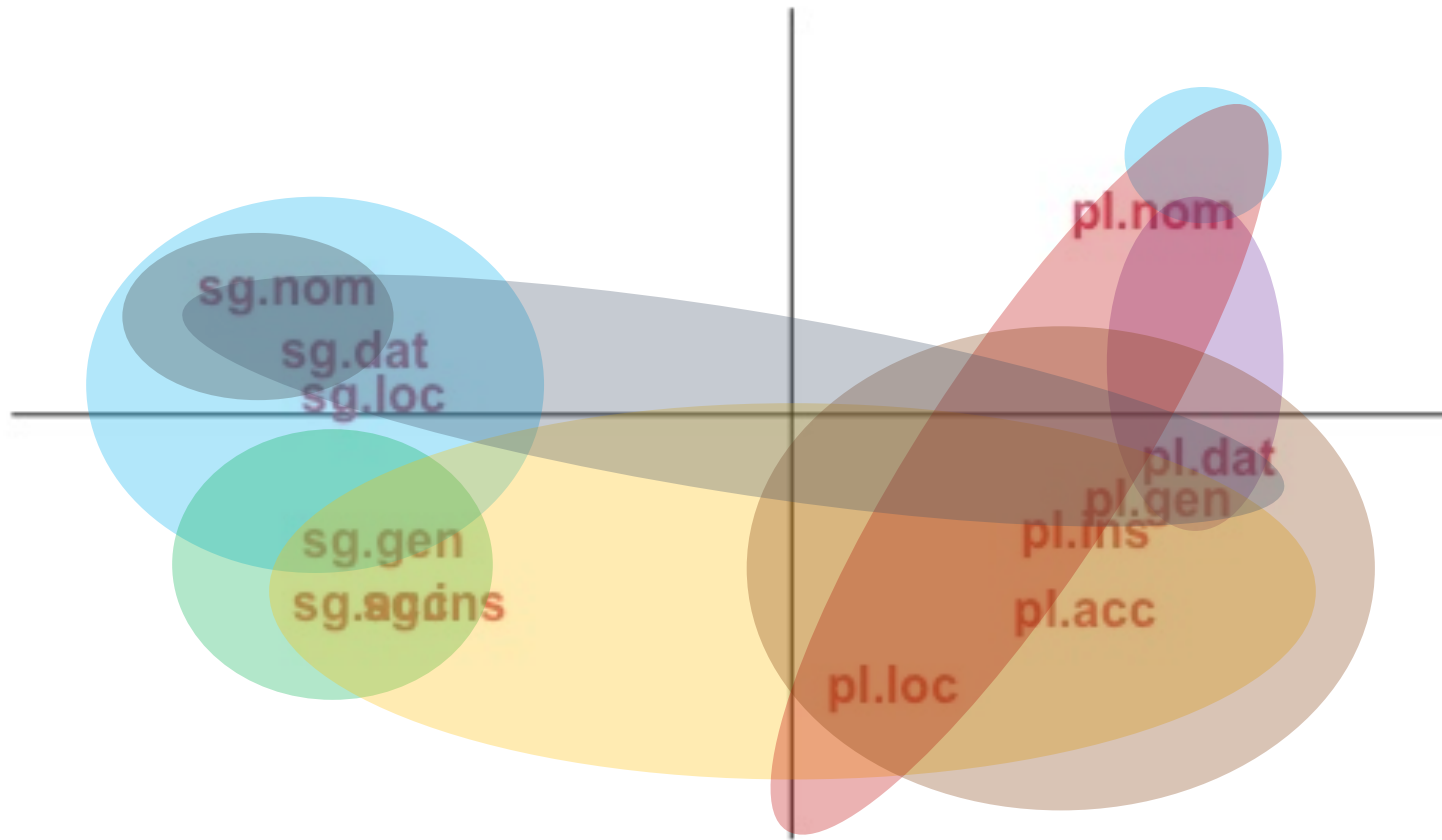


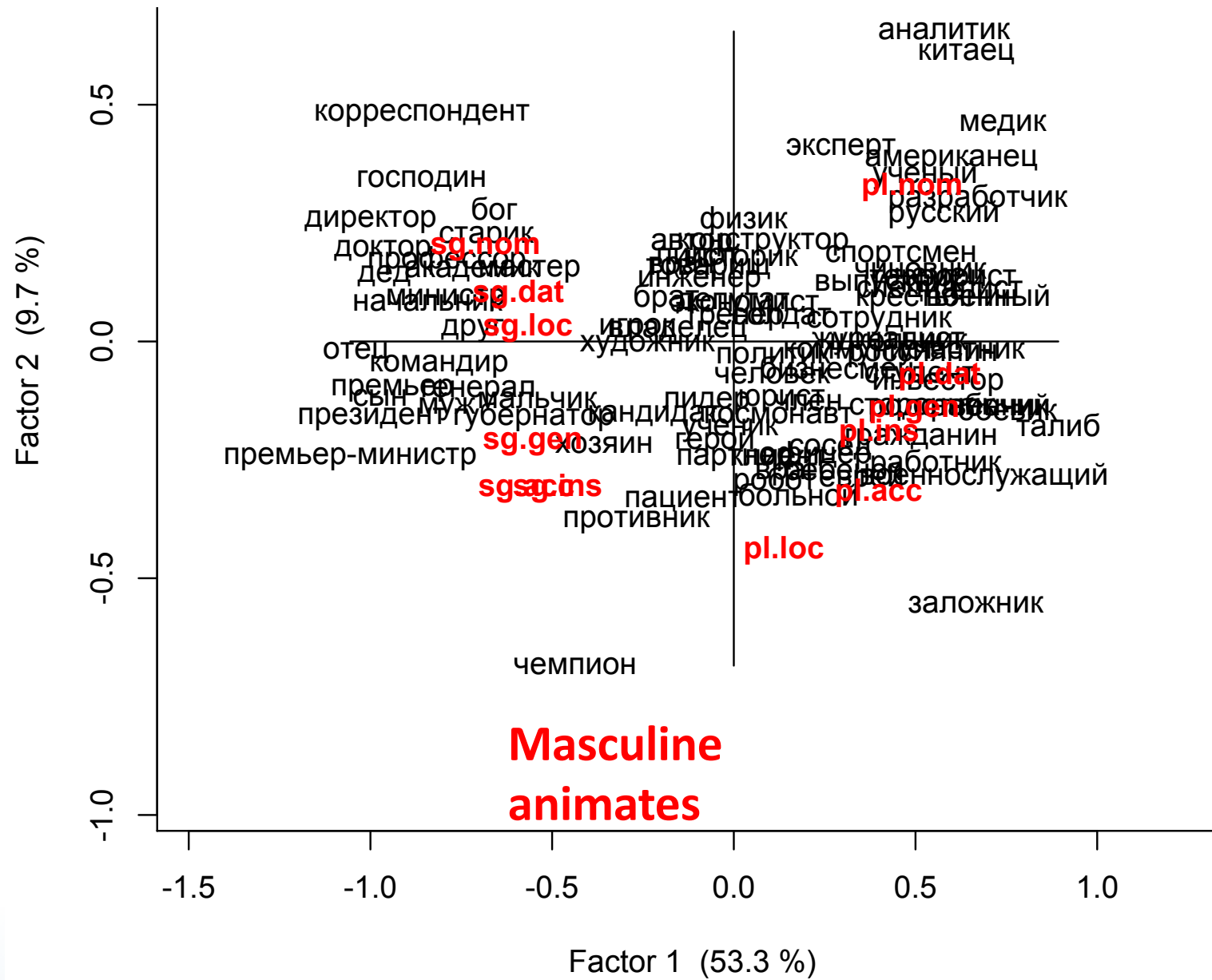












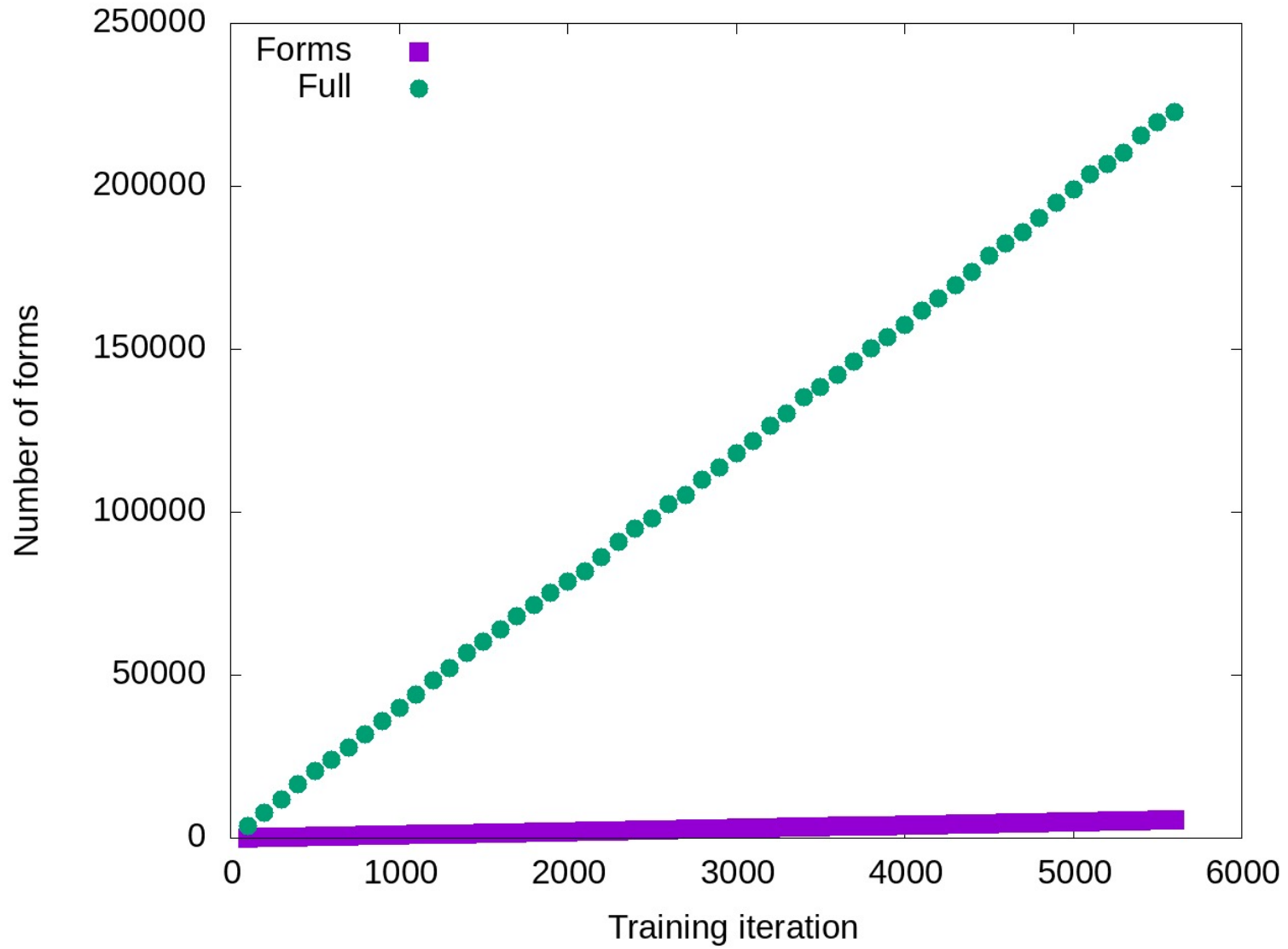
A machine-learning experiment

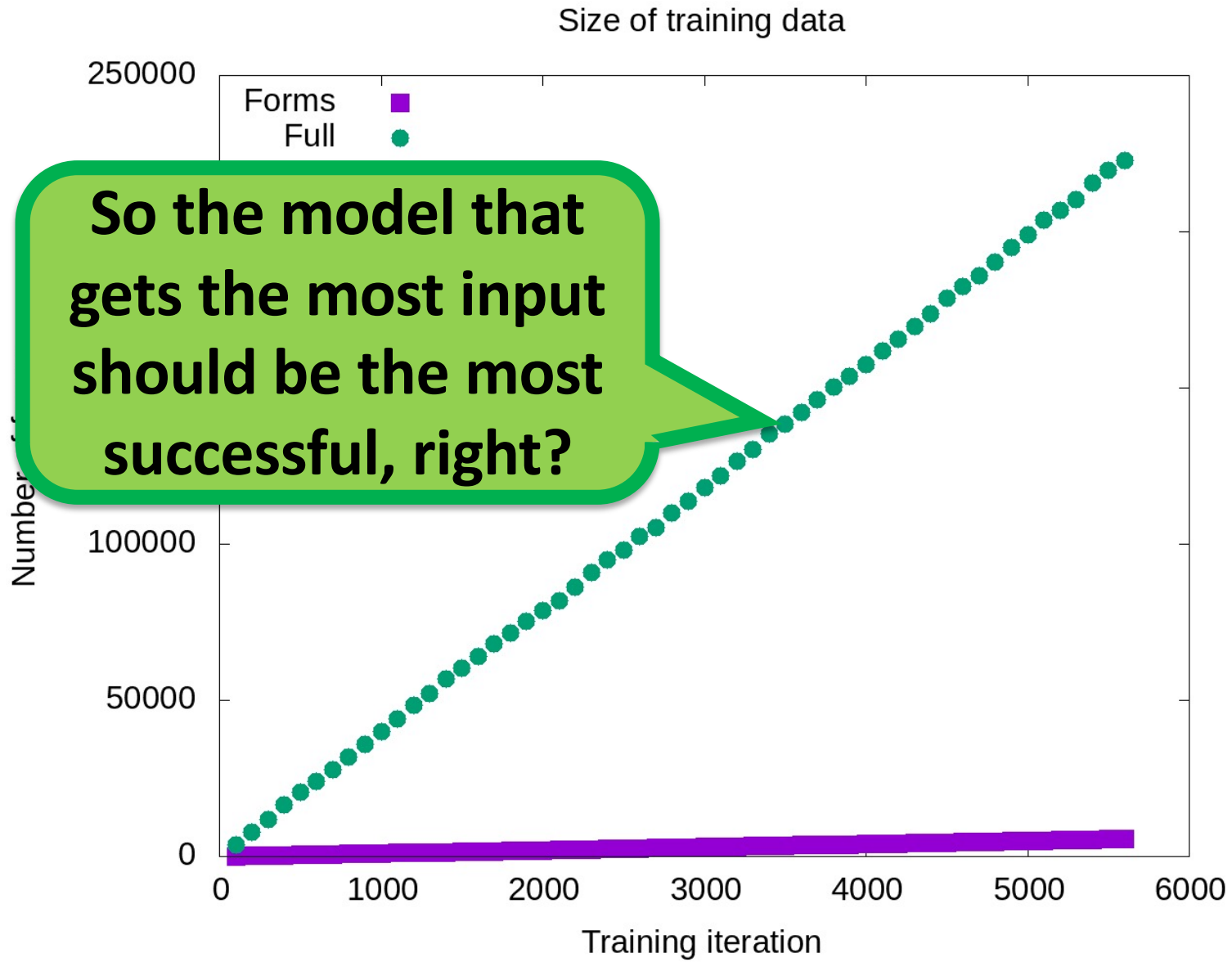
- Based on an ordered list of the most frequent forms for nouns, verbs, and adjectives in SynTagRus (1M gold standard)
- Machine learning: training, testing
 - Given the 100 most frequent forms, predict the next 100 most frequent forms
 - Given the 200 most frequent forms, predict the next 100 most frequent forms
 - Given the 300 most frequent forms, predict the next 100 most frequent forms
 - ... until 5400, when SynTagRus runs out of data
- Testing is always: Predict paradigm forms for 100 previously unseen words
- Two versions of experiment:
 - Training on entire paradigm, all forms
 - Training only on the single most frequent form

Data for training and testing from SynTagRus

Frequency & Form	Lemma	POS	Parse of form
1447 может	мочь	VERB	Aspect=Imp Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin Voice=Act
1286 года	год	NOUN	Animacy=Inan Case=Gen Gender=Masc Number=Sing
999 лет	год	NOUN	Animacy=Inan Case=Gen Gender=Masc Number=Plur
832 году	год	NOUN	Animacy=Inan Case=Loc Gender=Masc Number=Sing
813 время	время	NOUN	Animacy=Inan Case=Acc Gender=Neut Number=Sing
678 россия	россия	NOUN	Animacy=Inan Case=Gen Gender=Fem Number=Sing
571 могут	мочь	VERB	Aspect=Imp Mood=Ind Number=Plur Person=3 Tense=Pres VerbForm=Fin Voice=Act
571 люди	человек	NOUN	Animacy=Anim Case=Nom Gender=Masc Number=Plur
543 россия	россия	NOUN	Animacy=Inan Case=Loc Gender=Fem Number=Sing
436 является	являться	VERB	Aspect=Imp Mood=Ind Number=Sing Person=3 Tense=Pres VerbForm=Fin Voice=Act
416 случае	случай	NOUN	Animacy=Inan Case=Loc Gender=Masc Number=Sing
411 людей	человек	NOUN	Animacy=Anim Case=Gen Gender=Masc Number=Plur
403 страны	страна	NOUN	Animacy=Inan Case=Gen Gender=Fem Number=Sing
400 жизни	жизнь	NOUN	Animacy=Inan Case=Gen Gender=Fem Number=Sing

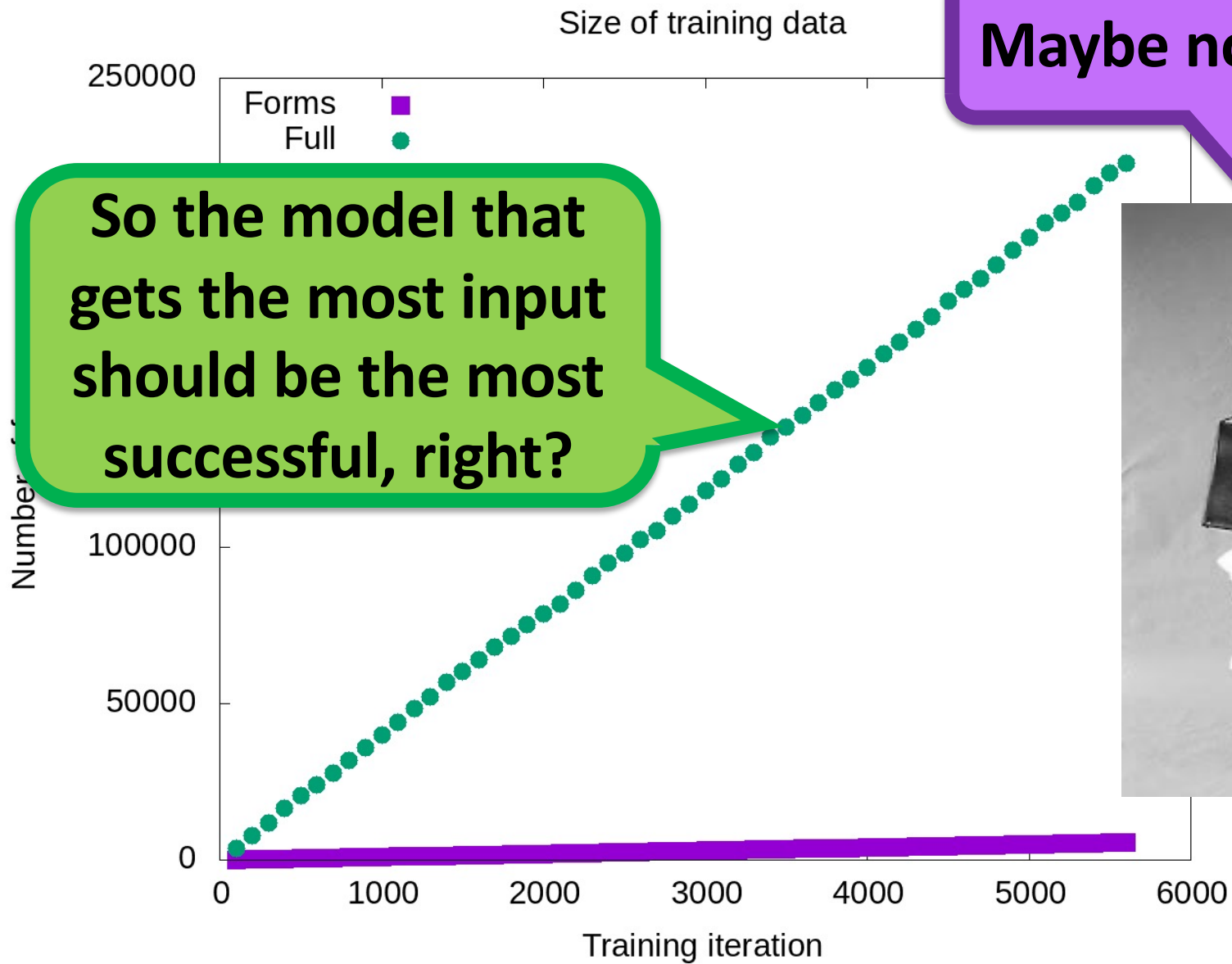
Size of training data





So the model that gets the most input should be the most successful, right?



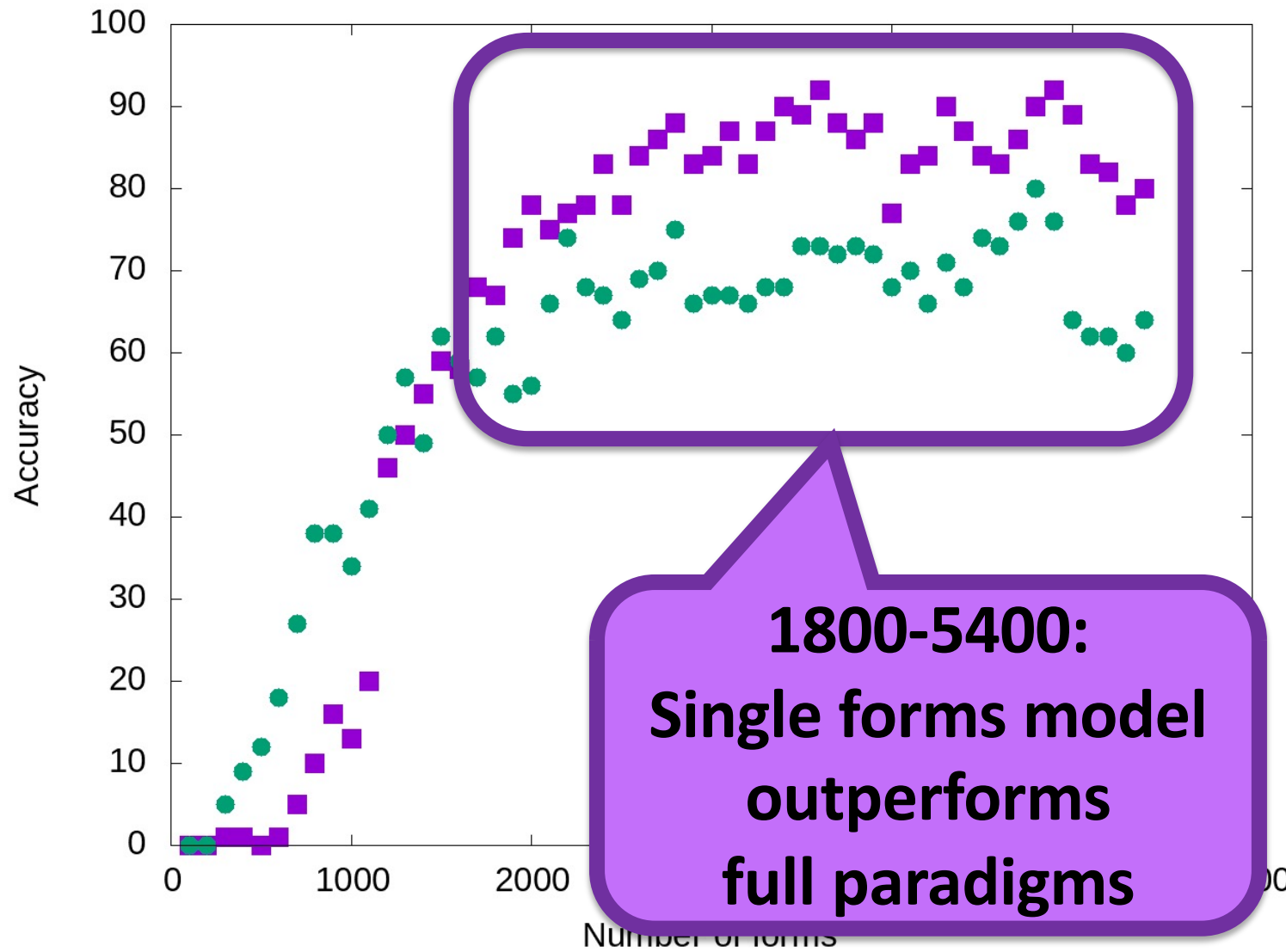


So the model that gets the most input should be the most successful, right?

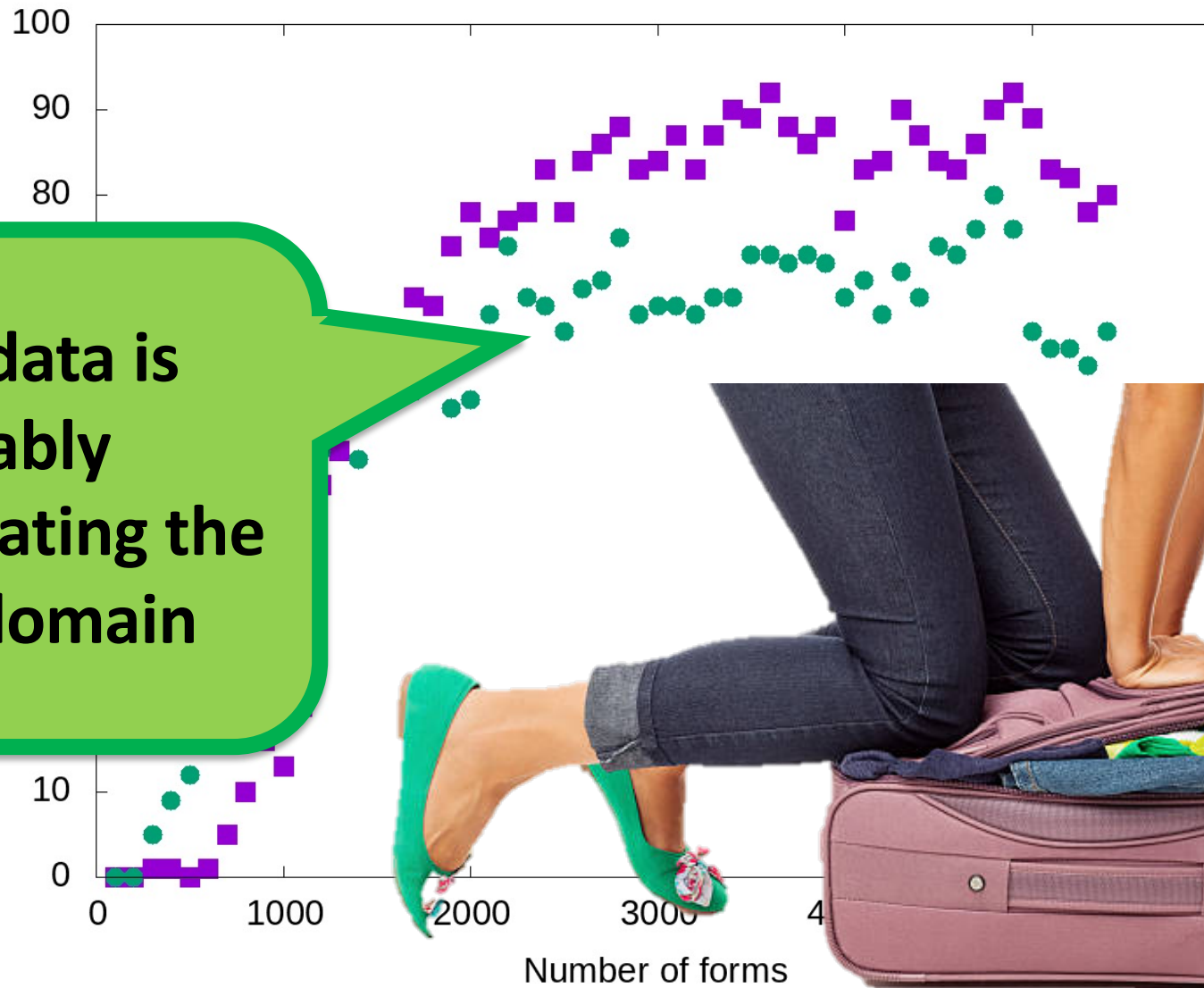
Maybe not...



Comparison of accuracy training on individual forms and full paradigms



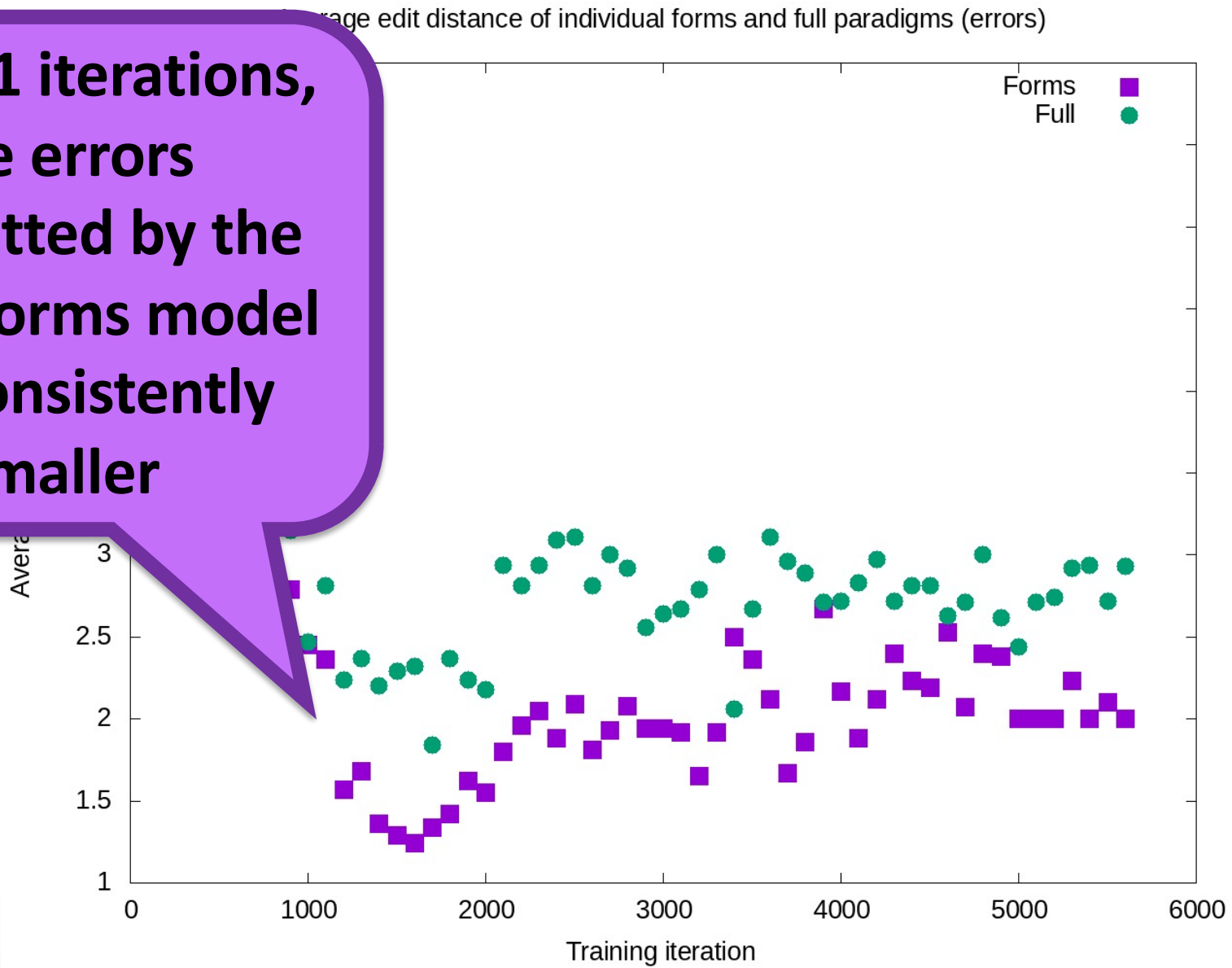
Comparison of accuracy training on individual forms and full paradigms



Excess data is probably overpopulating the search domain



**After 11 iterations,
the errors
committed by the
single forms model
are consistently
smaller**

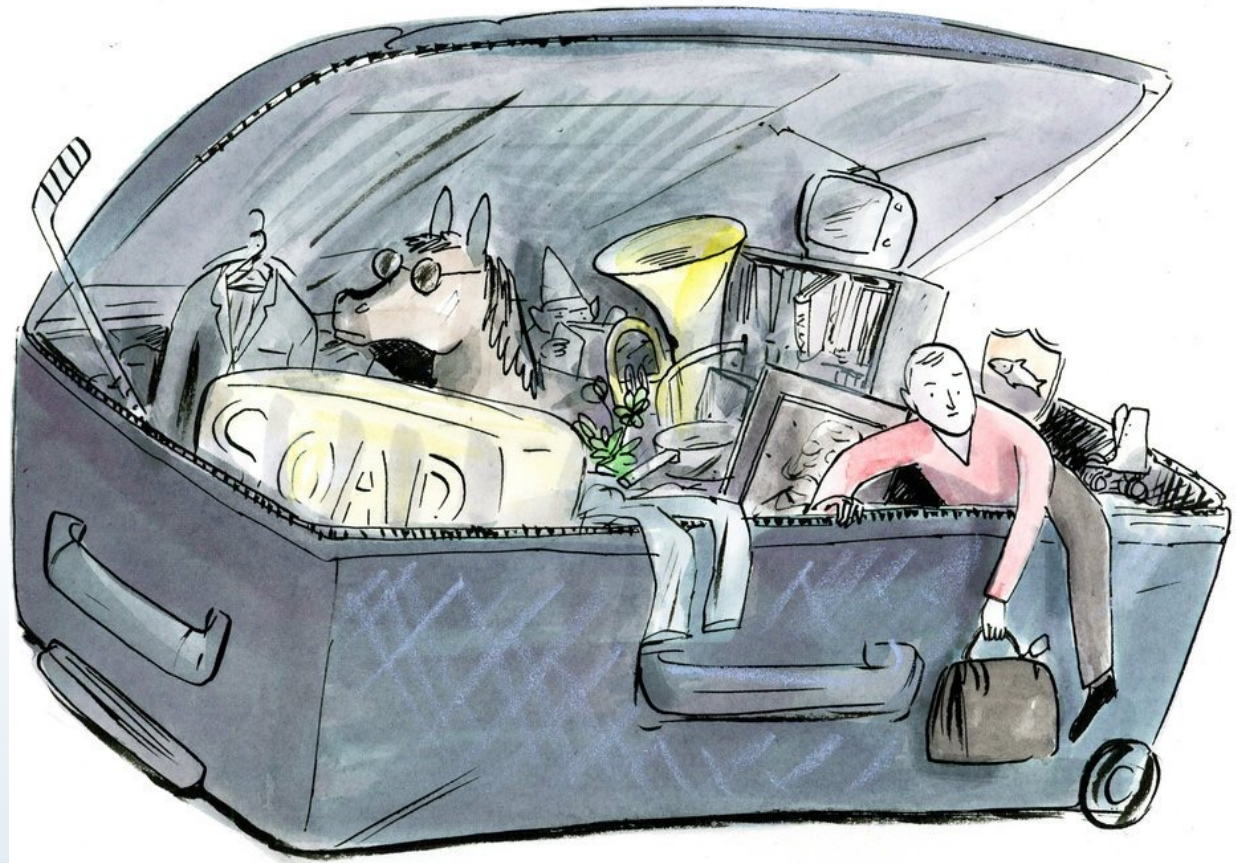


Machine-learning experiment: Summary

- Learning is potentially enhanced by focus **only on the most typical wordforms** attested for each lexeme: **accuracy increases** and **severity of errors decreases**
- This finding is **consistent with a usage-based cognitively plausible model**

How Can We Escape From Overstuffed Paradigms?

- Textbooks have always focused on certain forms and constructions
- Now we can do this in a scientific, consistent way



SMARTool = Strategic Mastery of Russian Tool

The SMARTool is available at: <http://uit-no.github.io/smartool/>

Free web resource for L2 learners of Russian that implements findings of corpus research and a learning simulation experiment to optimize the acquisition of Russian vocabulary and morphology

Partially funded by SIU [Center for Internationalization of Education], now known as DIKU [Norwegian Agency for International Cooperation and Quality Enhancement in Higher Education], in collaboration with the National Research University Higher School of Economics in Moscow



The SMARTool:

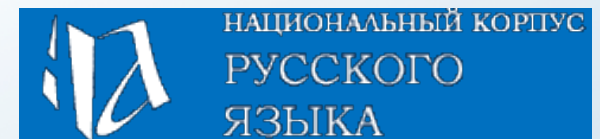
- **Inspired by research** on the distribution and simulated learning of Russian wordforms (**cognitively plausible**)
- Strategic focus on the **highest frequency wordforms and contexts** that motivate their use (**usage-based**)
- **Reduces the task** of learning a basic vocabulary of about 3,000 lexemes **by over 90%**
- Can be **continuously updated** and custom-tailored
- Potentially **portable to other languages** with rich inflectional morphology

Vocabulary Selection from 5 Textbooks and Лексический минимум; Balanced for Nouns, Verbs, Adjs (RNC ratio)

CEFR Level	ACTFL Equivalent	Russian Equivalent	SMARTool number of lexemes
A1 “Beginner”	Novice Low-Mid	ТЭУ Элементарный уровень	500
A2 “Elementary”	Novice High	ТБУ Базовый уровень	500
B1 “Intermediate”	Intermediate Low-Mid	ТРКИ-1 I Сертификационный уровень	1,000
B2 “Upper Intermediate”	Intermediate High-Advanced Low	ТРКИ-2	1,000

Typical Contexts Illustrated by Examples

- For each lexeme we identify 3 most common wordforms and **most typical grammatical constructions** and **lexical collocations**, and provide corpus-inspired **example sentences**
- Based on queries:
 - SynTagRus Corpus
 - the Russian National Corpus (<http://ruscorpora.ru>)
 - the Collocations Colligations Corpora (<http://cococo.cosyco.ru/>)
 - the Russian Constructicon (<https://constructicon.github.io/russian/>)



SMARTool

Search by topic

Search by analysis

Search by dictionary

List of abbreviations

About

Level

Topic

Show translation male voice female voice



1) Select a Level

2) Search by topic,
analysis, dictionary

Find the SMARTool here:

<https://smartool.github.io/smartool/>

A PLACE FOR US:



<https://dataverse.no/dataverse/trolling>



TROLLing

- is an international archive of linguistic data and statistical code
- is built on the Dataverse platform from Harvard University and complies with DataCite, the international standard for storing and citing research data
- is compliant with CLARIN, the EU research infrastructure for language-based resources
- assigns a permanent URL to each post
- uses metadata that ensures visibility and retrieval through international services
- is professionally managed by the University Library of Tromsø and an international steering committee



DataverseNO > TROLLing

Contact Share

Sign Up

Getting started with TROLLing



Search this dataverse...

Find

Advanced Search

Dataverses (0)

Datasets (111)

Files (2,998)

Publication Year

2014 (21)

2020 (19)

2016 (17)

2021 (14)

2017 (12)

More...

Author Name

1 to 10 of 111 Results

Sort ▾

Replication Data for: Eye Behavior during Syntactic Movement Evidence for Processing Approach to Persian Syntax

Jun 28, 2021

Alaee, Majid; Rasekh-Mahand, Mohammad; Tehrani-Doost, Mahdi, 2021, "Replication Data for: Eye Behavior during Syntactic Movement Evidence for Processing Approach to Persian Syntax", <https://doi.org/10.18710/TZBAOR>, DataverseNO, V1

[Dataset abstract:] The datasets include experimental material on an eye-tracking analysis in Persian syntax and how syntactic movement in this language is affected by weight factor. It also includes statistical datasets and the results of statistical analysis. In addition, some...

Replication Data for: Davvisámi earutkeahthes oamasteapmi

Jun 24, 2021

Janda, Laura A; Antonsen, Lene, 2020, "Replication Data for: Davvisámi earutkeahthes oamasteapmi", <https://doi.org/10.18710/QGXLOR>, DataverseNO, V2, UNE:6:TEN1nDnmWcpB9xsXi.lgcZg==.fileUNE1

Getting started with TROLLing

<http://site.uit.no/trolling/getting-started/>

- Promotional video
- Instructional videos
 - User guide
- TROLLing banner

Basic steps in TROLLing

- Create an account
 - Needed only for archiving -- you do not need an account to search or download data
 - This step is self-explanatory, but there is an instructional video
 - It may take a day or two for your account to be approved
- Create a study
 - Enter metadata, upload files in persistent formats, get DOI
- Search for a study
 - TROLLing terms of use

After you have uploaded your study...

- You submit it for approval and receive an acknowledgement
- It will be approved and released by an administrator and you will receive an email
- You will be able to edit your study later if needed and resubmit
- Previous versions of your study are archived, but only the latest version shows up in initial searches

Search for a study

- All of the cataloging information (metadata) is searchable, including:
 - author
 - affiliation
 - country / nation
 - date of production and distribution
 - keywords, e.g. language
 - topic classification
- Advanced search
 - possible to include and exclude combinations of items

Terms of use that users must agree to when downloading files

Affirm that:

- You will not use the materials that you downloaded to identify individuals or organizations
- You will not download or use materials that would be prohibited by law
- You will cite the data as stipulated in the Data Citation Information in any publications or reports that you make using the data