

The Biennial of Czech Linguistics, Data-based research in word formation
Charles University, Prague, September 20, 2024

Jurgis Pakerys (Vilnius University)

Virginijus Dadurkevičius (Vytautas Magnus University)

Agnė Navickaitė-Klišauskienė (Vilnius University)

How the measures of derivational productivity depend on lemmatization quality: the case of Lithuanian deverbal nouns

The project Derivational productivity of Lithuanian suffixed nouns in the Joint Corpus of Lithuanian has received funding from the Research Council of Lithuania (LMTLT), agreement No S-LIP-22-6

Outline

1. Derivational productivity measures
2. Our corpus and lemmatizer
3. Lemmatization and productivity measures
 - Fully automatic: agent nouns
 - Semi-automatic: agent nouns
 - Other categories
4. Conclusions
5. References



1. Derivational productivity measures

- **Realized** productivity = lemmas with affix X
- **Expanding** productivity = hapaxes with affix X or:
$$\frac{\text{hapaxes with affix } X}{\text{total hapaxes}}$$
- **Potential** productivity = $\frac{\text{hapaxes with affix } X}{\text{total frequency of lemmas with affix } X}$

Baayen (1992; 1993); overviews: Baayen (2009); Gaeta & Ricca (2015: 844–849), Dal & Namer (2016: 73–76), etc.

1. Derivational productivity measures

- Our claim: in the case of Lithuanian, **these measures are heavily dependent on the quality of lemmatization** (including automatic and manual stages)
- Similar observations:
 - German: Evert & Lüdeling (2001)
 - Italian: Gaeta & Ricca (2006)
 - French: Dal et al. (2008)
 - ...
 - In general: Baayen (2009: 207)

2. Our corpus and lemmatizer

- The Joint Corpus of Lithuanian, **1.3 billion tokens**
Dadurkevičius (2020a; 2020b), Dadurkevičius & Petrauskaitė (2020)

Subcorpus	Tokens
Lithuanian Internet content, 2014	779 M
Legal documents, Parliament of Lithuania, 2011	443 M
Balanced corpus of Vytautas Magnus University, 2008	113 M

- Lemmatized with Hunspell-based **lemmatizer: fixed dictionary + inflectional rules** (Dadurkevičius 2017)

3. Lemmatization & productivity measures

3.1. Fully automatic lemmatization

3.2. Additional semi-automatic lemmatization

- **Lithuanian deverbal suffixed nouns:** actions, agents, instruments, results/objects, places
- **Focus on the agent nouns** (3.1–3.2), brief notes on other categories (3.3)
- **Only realized and expanding productivity**

3.1 Fully automatic lemmatization

- Initial automatic lemmatization: **low type and hapax counts**
- For agent nouns, we still reviewed the data to determine how many non-transparent formations and accidentally included items have to be excluded
- For other categories, we omitted this stage and proceeded to additional semi-automatic lemmatization

Suffix	Before manual review		After manual review	
	V	V ₁	V	V ₁
<i>-toj-as</i> (m)	632	11	627	9
<i>-toj-a</i> (f)	547	11	543	9
<i>-ěj-as</i> (m)	251	2	235	2
<i>-ěj-a</i> (f)	258	3	201	2
<i>-ik-as</i> (m)	293	2	88	0
<i>-ik-ě</i> (f)	159	6	50	4

3.1 Fully automatic lemmatization

Agent nouns

V – types

V₁ – hapaxes

m – masculine

f – feminine

3.2 Additional semi-automatic lemmatization

- Tokens were automatically filtered according to pattern SUFFIX + ENDING and grouped into (potential) lemmas
- Lemma lists were manually reviewed, artificially constructed and derivationally non-transparent items were excluded
- Additional lemmatization significantly increased type and hapax counts

Suffix	Before manual review		After manual review	
	V	V ₁	V	V ₁
<i>-toj-as</i> (m)	3,305	822	2,456	532
<i>-toj-a</i> (f)	2,590	637	2,299	525
<i>-ěj-as</i> (m)	2,351	642	686	130
<i>-ěj-a</i> (f)	1,576	384	640	128
<i>-ik-as</i> (m)	911	189	254	65
<i>-ik-ě</i> (f)	857	274	92	27

3.2 Additional semi-automatic lemmatization

Agent nouns

V – types

V₁ – hapaxes

m – masculine

f – feminine

3.2 Additional semi-automatic lemmatization

- **Type and hapax counts** were found to be somewhat **inflated due to homographic forms** of masc./fem. nouns
- Morphologically ambiguous tokens yielded two lemmas due to the lack of morphological disambiguation stage during the lemmatization
- **We manually disambiguated the hapaxes**
- This step significantly reduced hapax counts of fem. formations in *-toj-a* and *-éj-a* due to high number of their shared homographic forms

Suffix	Before manual review		After manual review		
	V	V ₁	V	V ₁ (before m/f disambiguation)	V ₁ (after m/f disambiguation)
<i>-toj-as</i> (m)	3,305	822	2,456	532	464
<i>-toj-a</i> (f)	2,590	637	2,299	525	65
<i>-èj-as</i> (m)	2,351	642	686	130	112
<i>-èj-a</i> (f)	1,576	384	640	128	24
<i>-ik-as</i> (m)	911	189	254	65	61
<i>-ik-è</i> (f)	857	274	92	27	20

3.2 Additional semi-automatic lemmatization

Agent nouns

V – types

V₁ – hapaxes

m – masculine

f – feminine

3.3 Other categories

Suffix	Initial automatic lemmatization		Additional semi-automatic lemmatization			
	No manual review		Before manual review		After manual review	
	V	V ₁	V	V ₁	V	V ₁
<i>-im-as(-is)</i>	9,634	610	17,134	3,591	15,274	2,855
<i>-ym-as(-is)</i>	1,770	234	2,101	418	1,943	351
<i>-acij-a</i>	795	15	1,145	113	966	69
<i>-es-ys</i>	66	0	350	105	154	27

Action nominals

V – types

V₁ – hapaxes

Suffix	Initial automatic lemmatization		Additional semi-automatic lemmatization			
	No manual review		Before manual review		After manual review (cases of agent/instrument/ place nouns still not disambiguated)	
	V	V ₁	V	V ₁	V	V ₁
<i>-tuv-as</i>	192	3	592	127	488	87
<i>-ikl-is</i>	134	4	332	67	293	54
<i>-atori-us</i>	239	6	468	69	208	16
<i>-tuv-é</i>	37	0	248	51	123	20

3.3 Other categories

Instrument nouns
V – types
V₁ – hapaxes

Agent/instrument categories in 3.2 were also not disambiguated

3.3 Other categories

Suffix	Initial automatic lemmatization		Additional semi-automatic lemmatization			
	No manual review		Before manual review		After manual review	
	V	V ₁	V	V ₁	V	V ₁
<i>-in-ys</i>	193	1	837	143	230	16
<i>-al-as</i>	260	5	618	126	124	31
<i>-at-as</i>	466	16	2,222	514	72	5

Object/result nouns

V – types

V₁ – hapaxes

3.3 Other categories

Suffix	Initial automatic lemmatization		Additional semi-automatic lemmatization			
	No manual review		Before manual review		After manual review (and disambiguation of place/instrument/celebration)	
	V	V ₁	V	V ₁	V	V ₁
<i>-ykl-a</i>	89	0	218	52	189	41
<i>-tuv-é</i>	37	0	248	51	33	7

Place nouns

V – types

V₁ – hapaxes

4. Conclusions

- The quality of lemmatization significantly influences the estimates of derivational productivity
- Lemmatization parameters are very important and should be listed in all corpus descriptions:
 - Built-in dictionaries, word-guessing modules, capabilities of solving morphological ambiguity, machine learning techniques, etc.

Thank you!

jurgis.pakerys@flf.vu.lt



5. References

Baayen, R. Harald (1992). Quantitative Aspects of Morphological Productivity. In: Booij, Geert E. & van Marle, Jaap (eds.), *Yearbook of Morphology 1991*. Dordrecht: Kluwer Academic Publishers, 109–149 (https://doi.org/10.1007/978-94-011-2516-1_8).

Baayen, R. Harald (1993). On Frequency, Transparency, and Productivity. In: Booij, Geert E. & van Marle, Jaap (eds.), *Yearbook of Morphology 1992*. Dordrecht: Kluwer Academic Publishers, 181–208 (https://doi.org/10.1007/978-94-017-3710-4_7).

Baayen, R. Harald (2009). Corpus linguistics in morphology: Morphological productivity. In: Lüdeling, Anke & Kytö, Merja (eds.), *Corpus Linguistics: An International Handbook*, vol. 2. Berlin, New York: Mouton de Gruyter, 899–919 (<https://doi.org/10.1515/9783110213881.2.899>).

Dadurkevičius, Virginijus (2017). Lietuvių kalbos morfologija atvirojo kodo „Hunspell“ platformoje, *Bendrinė kalba* 90, 1–17 (<https://journals.iki.lt/bendrinekalba/article/view/156>).

Dadurkevičius, Virginijus (2020a). Wordlist of Lemmas from the Joint Corpus of Lithuanian. *CLARIN-LT digital library in the Republic of Lithuania* (<http://hdl.handle.net/20.500.11821/41>).

Dadurkevičius, Virginijus (2020b). Assessment data of the Dictionary of Modern Lithuanian versus Joint Corpora, *CLARIN-LT digital library in the Republic of Lithuania* (<https://clarin.vdu.lt/xmlui/handle/20.500.11821/36>).

Dadurkevičius, Virginijus & Petrauskaitė, Rūta (2020). Corpus-based methods for assessment of traditional dictionaries. In: Utkā, Andrius, Vaičėnonienė, Jurgita, Kovalevskaitė, Jolanta & Kalinauskaitė, Danguolė (eds.), *Human Language Technologies—The Baltic Perspective. Frontiers in Artificial Intelligence and Applications*, vol. 328. Amsterdam: IOS Press, 123–126 (<https://doi.org/10.3233/FAIA200613>).

Dal, Georgette, Fradin, Bernard, Grabar, Natalia, Namer, Fiammetta, Lignon, Stéphanie, Plancq, Clément, Zweigenbaum, Pierre & Yvon, François (2008). Quelques préalables au calcul de la productivité des règles constructionnelles et premiers résultats. In: Durand, Jacques, Habert, Benoît & Laks, Bernard (eds.), *Actes du premier Congrès mondial de linguistique française, Paris, 9–12 juillet 2008*. Paris: Institut de Linguistique Française, 1587–1599 (<https://doi.org/10.1051/cmlf08184>).

Dal, Georgette & Namer, Fiammetta (2016). Productivity. In Hippisley, Andrew & Stump, Gregory (eds.), *The Cambridge Handbook of Morphology*. Cambridge: Cambridge University Press, 70–90 (<https://doi.org/10.1017/9781139814720.004>).

Evert, Stefan & Lüdeling, Anke (2001). Measuring morphological productivity: Is automatic preprocessing sufficient? In: Rayson, Paul, Wilson, Andrew, McEnery, Tony, Hardie, Andrew & Khoja, Shereen (eds.), *Proceedings of the Corpus Linguistics 2001 Conference*. Lancaster: Lancaster University, 167–175.

Gaeta, Livio & Ricca, Davide (2006), Productivity in Italian word formation: a variable-corpus approach. *Linguistics* 44(1), 57–89 (<https://doi.org/10.1515/LING.2006.003>).

Gaeta, Livio & Ricca, Davide (2015). Productivity. In: Müller, Peter O., Ohnheiser, Ingeborg, Olsen, Susan & Rainer, Franz (eds.), *Word-Formation: An International Handbook of the Languages of Europe*, vol. 2. Berlin/Boston: De Gruyter Mouton, 842–858 (<https://doi.org/10.1515/9783110246278-003>).