



# Community proofreading as a tool for community engagement

A quantitative analysis

Sebastian Nordhoff

June 3, 2019

FU Berlin

- › Open Access is mainly concerned with reading
- › Open Publishing is concerned with making all aspects of publishing open (Rob Cartolano)
  - › Open source platforms
  - › Open formats
  - › Open protocols
  - › Open bookkeeping
  - › Open peer review
  - › Community proofreading

- › one research can adopt different roles
  - › author, reviewer, reader, ...
- › junior researchers are more often readers
- › senior researchers take on the other roles as well
- › complex ecosystem
- › community-based publishing tries to integrate researchers at all levels

- › outsourced work-for-hire
- › for a fee
- › one proofreader
- › specialist in style and guidelines
- › might have some training in linguistics
- › normally no specialist knowledge of the particular subfield

- › crowdsourced to the community
- › voluntary work
- › many proofreaders, often junior
- › very often specialists in the particular subfield
- › intrinsic interest
- › less acquaintance with style and guidelines

- › Open Access publisher in linguistics
- › 100+ books since 2014
- › 350 community proofreaders

<p><b>Natural causes of language: Frames, biases, and cultural transmission</b> N. Enfield (Author)</p>	<p><b>A topology of marked languages</b> Corina Fiebert (Author)</p>	<p><b>The Talking Heads experiment: Origins of words and meanings</b> L. L. Lewis (Author)</p>	<p><b>Grammaticalization in the North: Noun phrase morphemes in Scandinavian vernaculars</b> Ole Dahl (Author)</p>	<p><b>A grammar of Fula (Sami)</b> Jonas Wihor (Author)</p>	<p><b>Phonetic detail in Neogloss</b> Francesco Cengini (Author)</p>
<p><b>Linguistic variation, identity construction and cognition</b> Karin K. Dräger (Author)</p>	<p><b>Language strategies for the domain of labor</b> Jonas Wihor (Author)</p>	<p><b>A grammar of Fula (Sami)</b> Ole Dahl (Author)</p>	<p><b>Thoughts on grammaticalization</b> Christoph Lehmann (Author)</p>	<p><b>How mobile robots can self-organize a vocabulary</b> Paul Fiegler (Author)</p>	<p><b>A grammar of Mawé</b> Lisa Berglund (Author)</p>
<p><b>Roots of language</b> Derek Bickerton (Author)</p>	<p><b>The empirical base of linguistics: Grammatical judgments and linguistic methodology</b> Carsten T. Schuler (Author)</p>	<p><b>A grammar of Fula (Sami)</b> Ole Dahl (Author)</p>	<p><b>New directions in corpus-based translation studies</b> Christy Farnsworth, Frederic Zeman (Volume Editor)</p>	<p><b>Grammatical theory: From transformational grammar to construction-based approaches</b> Stefan Müller (Author)</p>	<p><b>Advances in the study of Slavic languages and linguistics</b> Catherine Ruy, Björn Petersen (Volume Editor)</p>
<p><b>The evolution of grounded spatial language</b> Michael Spranger (Author)</p>	<p><b>The evolution of case grammar</b> Bernard Trapp (Author)</p>	<p><b>The future of dialects: Selected papers from Methods in Dialectology IV</b> Marie-Hélène Côté, Bernice Knoch, John Nerbonne (Volume Editor)</p>	<p><b>Syntax und Valenz: Der Modellierung von Verbstrukturalen und affektiven Strukturen im Baumadjunktogramm</b> Thomas Lohr (Author)</p>	<p><b>Adjective attribution</b> Michael Keller (Author)</p>	<p><b>A grammar of Kanyo</b> Quinn (Author)</p>
<p><b>Sprachliche Sozialisation: Jüdischkeit in der deutschsprachigen Literatur (18.-20. Jahrhundert)</b> Lisa Schäfer (Author)</p>	<p><b>Forthcoming: Language technologies for environmental change</b> Gregory Bateson, David G. Bateson, John S. Bateson (Volume Editor)</p>	<p><b>Eyetracking and Applied Linguistics</b> Silvia Hansen-Schiera, Sandra Gracia (Volume Editor)</p>	<p><b>Forthcoming: Grammatical structures in comparative phrase structure analysis</b> Gabriela Wille (Author)</p>	<p><b>Einführung in die grammatische Beschreibung der Deutschen: Zweite, überarbeitete Auflage</b> Roland Schäfer (Author)</p>	<p><b>New perspectives on cultural coherence: Implications for translation</b> Karin M. M. D. D. (Volume Editor)</p>
<p><b>Time in Yenching: A lexical study and morphological analysis</b> Alexa Michael (Author)</p>	<p><b>A grammar of Kanyo</b> Quinn (Author)</p>	<p><b>Forthcoming: Attribution construction in North-Eastern Neo-Aramaic</b> Jonas Wihor (Author)</p>	<p><b>Forthcoming: Empirical modelling of translation and interpreting</b> Silvia Hansen-Schiera, Oliver Gals, Sascha Hoffmann, Bernd Meyer (Volume Editor)</p>	<p><b>Forthcoming: Reflections on Linguistics</b> Gerd B. B. (Author)</p>	<p><b>Diversity in African languages: Selected papers from the 4th Conference on African Linguistics</b> Doris L. Payne, Sara P. Payne (Volume Editor)</p>
<p><b>Annotation, exploitation and evaluation of parallel corpora: TC3</b> Silvia Hansen-Schiera, Sandra Gracia (Volume Editor)</p>	<p><b>Forthcoming: Beiträge zur deutschen Grammatik</b> Tobias H. H. (Author), Stefan Müller, Frank Richter, Margarete (Volume Editor)</p>	<p><b>The Bilingual Dictionary and grammar sketch</b> Tobias H. H. (Author)</p>	<p><b>A grammar of Fula (Sami)</b> Jonas Wihor (Author)</p>	<p><b>Forthcoming: Unity and diversity in grammaticalization scenarios</b> Walter Stroh, Achim Machobane (Volume Editor)</p>	<p><b>Tonal placement in Tachibana</b> Tachibana (Author)</p>
<p><b>Forthcoming: Order and structure in syntax II: Subordinate and argument structure</b> Laura A. Bailey, Michelle Sheehan (Volume Editor)</p>	<p><b>Forthcoming: Order and structure in syntax I: Word order and syntactic structure</b> Laura A. Bailey, Michelle Sheehan (Volume Editor)</p>	<p><b>Forthcoming: A equívoco da língua materna em português: Questões gerais e específicas da Português</b> Marta João Freitas, Ana Lúcia Santos (Volume Editor)</p>	<p><b>Forthcoming: Further investigations in the syntax of infinitives</b> Cecilia T. (Volume Editor)</p>	<p><b>Dependencies in language: On the central argument of linguistic systems</b> N. J. Enfield (Volume Editor)</p>	<p><b>On looking towards (and beyond) the future</b> Doris L. Payne, Sara P. Payne (Volume Editor)</p>

- › proofreading queue with a new title every 2 weeks
- › title is announced on Monday
- › community members can volunteer and claim a chapter
- › chapters are assigned on Wednesday
- › 4 weeks time for proofreading
- › proofreading is done on Paperhive



Document

Discussions 317

Activity

which would yield only default agreement on the RC predicate, contrary to fact.

## 2.2 The Genitive of Quantification phenomenon

The Genitive of Quantification phenomenon has been described to a large extent for Slavic languages in Bošković (2006); Franks (1994; 2002); Przepiórkowski (2004); Rutkowski (2002); and Willim (2003), to name but a few. In Polish, genitive case marking is forced on a noun which is modified by a higher numeral or a lower virile numeral, as well as by certain quantifiers such as *wiele* 'many', *kilka* 'a few', *para* 'a couple of', etc. Such numeral phrases do not induce subject-verb agreement in main clauses, as can be seen in (17), in which the verb obligatorily appears in the 3SG neuter form, regardless of the grammatical gender of the noun.

- (17) a. Siedmiu mężczyzn weszło/\*weszli do domu.  
seven.ACC men.GEN,VIR entered.3SG,NEUT/3PL,VIR into house  
'Seven men entered the house.'
- b. Siedem kobiet weszło/\*weszły do domu.  
seven.ACC women.GEN, NON-VIR entered.3SG,NEUT/\*3PL,NON-VIR into house  
'Seven women entered the house.'

The analysis of Polish GoQ structures proposed in Witkoś & Dziubała-Szrejmbrowska (2016) follows the idea that probing for phi-features is possible for T

x



Alexandr Rosen · 18 days ago

**N instead of NEUT**

according to

<https://www.eva.mpg.de/lingua/resources/glossing-rules.php>

Show more



Alexandr Rosen · 18 days ago

**GEN.VIR**

period as the separator

according to

<https://www.eva.mpg.de/lingua/resources/glossing-rules.php>

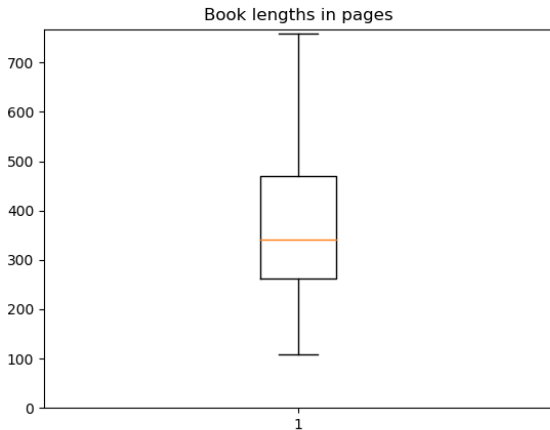
Show more

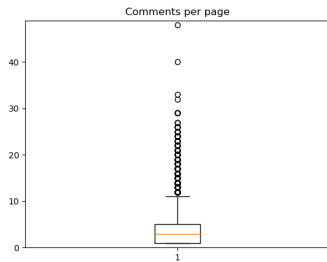
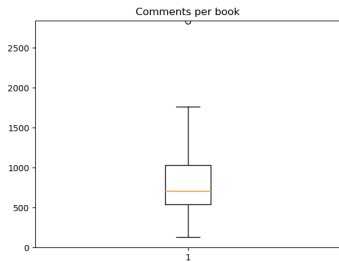


› Westedt analysed a sample of comments on Paperhive for her BA thesis.

Category	Percentage
Style	21.00
Lexical choice	20.73
Punctuation	11.81
Grammar	11.55
References	9.71
Syntax	7.80
Spelling	7.30
<b>Content</b>	<b>6.56</b>
Miscellanea	3.41

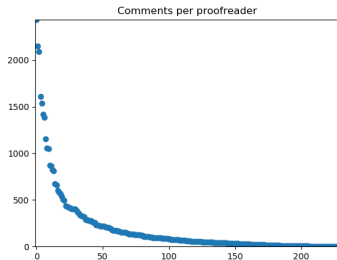
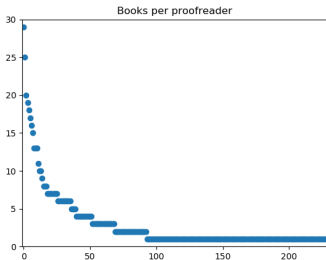
- › 52 books from late 2016 to late 2018
- › comments were harvested from Paperhive and put into a database
- › 19 004 pages
- › 43 370 comments



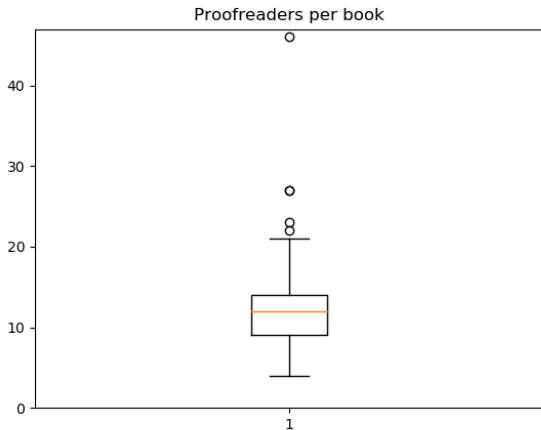


The highest number of comments on one page is found in Theory and description in African Linguistics on page 122 (48 comments).

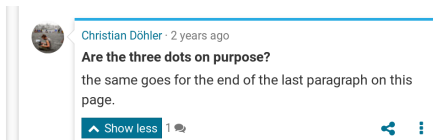
# Productivity of proofreaders



228 different accounts have participated in commenting.



course about the Jewish Nerwa texts, where I  
rstood it as a recommendation to further my  
this...

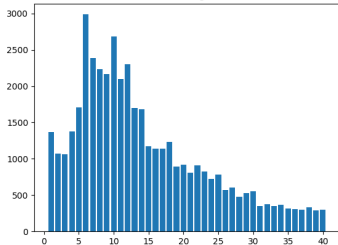


- › A PaperHive comment has a succinct title (<40 characters)
- › optional body, with more elaborate information

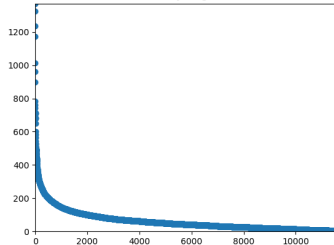


# Title length and body length

Title length



Body length

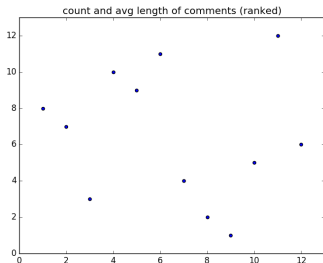


1. **Proofreaders fall into two types. Type 1 will focus on small details; type 2 will focus on the big picture.**
2. **Proofreading will diminish as the proofreader moves along. Comments will become shorter due to fatigue, i.e. average comment length will go down due to repetition of previous remarks as “see above”.**

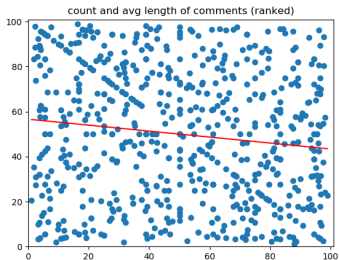
- › Type 1: many comments but short (“comma missing”)
- › Type 2: few comments, but longer, in-depth

› For every book

- › rank all participating proofreaders by amount of comments
- › rank all participating proofreaders by average length of comments
- › plot the two against each other



- > 12 proofreaders participated
- > their respective ranks are given by the dots.
  - > e.g. #3 in one rank is also #3 in the other, but #1 on one is #8 in the other
- > data from one book insufficient



- > Ranks are normalized to centiles
- > best fit given by red line
- > indeed a weak negative correlation

} Hypothesis #1 is confirmed

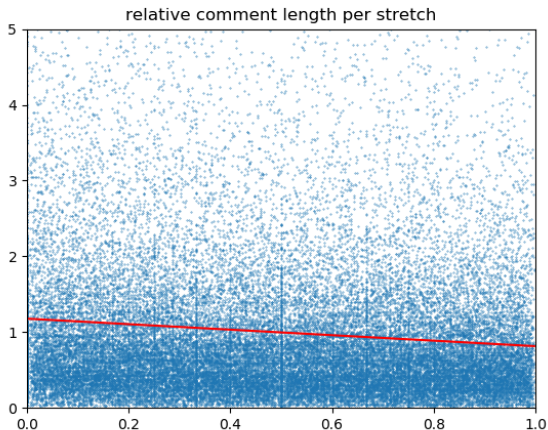
- } proofreaders with more comments have shorter comments
- } proofreaders with longer comments comment less

**Hypothesis 2 : Proofreading will diminish as the proofreader moves along. Comments will become shorter due to fatigue, i.e. average comment length will go down due to repetition of previous remarks as “see above”.**



- › for every book
  - for every proofreader
    - for every comment
      - › compute relative length (e.g. 0.67 of the average)
      - › compute relative position (front, middle, back)
      - › store the tuple (relative position, relative length)
      - › A dot at (0.5, 5) means that there was a comment in the middle of the relevant stretch whose length was 5 times the average comment length.
- › the relative position can be pegged to the linear order of comments, or to the pages

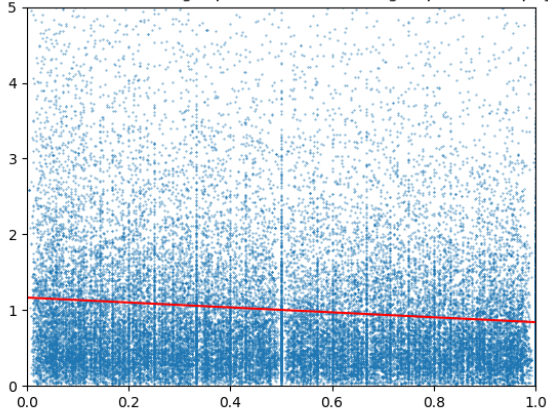
## Plot for Hypothesis #2 based on linear order



## Plot for Hypothesis #2 based on page position

| Hypothesis evaluation  
| Proofreader fatigue

relative comment length per stretch according to position of page



### } Hypothesis is confirmed

- } the later in the document a comment is, the shorter it will be
- } the first comment will be about 110% of the average, while the last one will be 90% of the average.
- } effect not very strong, but discernible

- › Main aim: methodological
- › Proofreading comments are a by-product of open publishing
  - › In traditional publishing models, these data would not be available
- › Once the documents, processes, and formats are opened up, novel research questions can emerge which would not have been possible under a closed setup.
- › Implications for psychology of reading for instance.

- › There are 908 people with the role “author” at LangSci Press
- › There are 228 proofreaders
- › 27 researchers have taken up both roles
  - › 16 started as authors, and became proofreaders later
  - › 11 started as proofreaders, and became authors later
  - › Movement between the author pool and the proofreader pool in both directions.

- › Community proofreading is a novel way of engaging the community
- › only possible for Open Access publications
- › workable implementation with 50+ books and 200+ researchers
- › can compare to traditional proofreading
- › by-product data can be used for novel research questions
  - › proofreader typology
  - › proofreader fatigue
- › flow back and forth between the group of authors and the group of proofreaders
- › healthy ecosystem
- › researchers from different backgrounds at different stages of their career contribute their respective expertises to creating and improving manuscripts.

- › What other questions could be addressed with that data?
- › Which other disciplines might be interested?





### Gold proofreaders

> **Andreas Hölzl** ([view profile](#))



34/99

> **Jeroen van de Weijer** ([view profile](#))



31/99

> **Eitan Grossman** ([view profile](#))

29/99

> **Jean Nitzke**

25/98

> **Christian Döhler** ([view profile](#))

19/97

> **Martin Haspelmath** ([view profile](#))

19/97

> **Ahmet Bilal Özdemir** ([view profile](#))

19/97

> **Ikmi Nur Oktavianti** ([view profile](#))

15/96

> **Brett Reynolds** ([view profile](#))

14/96