



Royal Netherlands Academy of Arts and Sciences (KNAW) KONINKLIJKE NEDERLANDSE AKADEMIE VAN WETENSCHAPPEN

The desirability of a corpus of online book responses

Boot, P.

published in

Proceedings of the Workshop on Computational Linguistics for Literature
2013

document version

Publisher's PDF, also known as Version of record

document license

CC BY

[Link to publication in KNAW Research Portal](#)

citation for published version (APA)

Boot, P. (2013). The desirability of a corpus of online book responses. In *Proceedings of the Workshop on Computational Linguistics for Literature* (pp. 32-40). Association for Computational Linguistics (ACL).
<http://aclweb.org/anthology/W/W13/W13-1405.pdf>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the KNAW public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain.
- You may freely distribute the URL identifying the publication in the KNAW public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

E-mail address:

pure@knaw.nl

The desirability of a corpus of online book responses

Peter Boot

Huygens ING

PO Box 90754

2509 HT The Hague

The Netherlands

`peter.boot@huygens.knaw.nl`

Abstract

This position paper argues the need for a comprehensive corpus of online book responses. Responses to books (in traditional reviews, book blogs, on booksellers' sites, etc.) are important for understanding how readers understand literature and how literary works become popular. A sufficiently large, varied and representative corpus of online responses to books will facilitate research into these processes. This corpus should include context information about the responses and should remain open to additional material. Based on a pilot study for the creation of a corpus of Dutch online book response, the paper shows how linguistic tools can find differences in word usage between responses from various sites. They can also reveal response type by clustering responses based on usage of either words or their POS-tags, and can show the sentiments expressed in the responses. LSA-based similarity between book fragments and response may be able to reveal the book fragments that most affected readers. The paper argues that a corpus of book responses can be an important instrument for research into reading behavior, reader response, book reviewing and literary appreciation.

1 Introduction

The literary system does not consist of authors and works alone. It includes readers (or listeners) and their responses to literary works. Research into reception is an important subfield of literary studies (e.g. Goldstein and Machor, 2008). Shared

attention to stories may have evolved as way of learning to understand others and to increase bonding (Boyd, 2009). Discussing literature may thus be something that we are wired to do, and that we do indeed wherever possible: today on Amazon, on weblogs, and on Twitter, and in earlier days in newspapers and letters. These responses to books are important both as documentation of the ways literary works are read and understood, and because they help determine works' short- and long-term success.

This position paper argues that what we need, therefore, is a large and representative corpus of book responses. 'Book response' in this paper includes any opinion that responds to a book, i.e. traditional book reviews, book-based discussion, opinions given on booksellers' sites, on Twitter, thoughtful blog posts, and the like. The word 'books' here is meant to refer to all genres, including literature as well as more popular genres such as fantasy, thrillers, comics, etc. Section 2 of the paper discusses the importance and research potential of book responses. Section 3 reviews related research. In section 4, I outline the properties that this corpus should have. Section 5 describes a Dutch pilot corpus and shows some aspects of this corpus that lend themselves to analysis with linguistic and stylometric tools. Section 6 presents conclusions and directions for future work.

The author of this paper is not a computational linguist, but has a background in literary studies and digital humanities. The intention is to create a dialogue between literary studies and computational linguistics about fruitful ways to investigate book responses, their relations to the books they

respond to and their effects on short-term or long-term appreciation.

2 Book responses and their importance

Evaluating books and talking about our response is a very natural thing to do (Van Peer, 2008). In a professionalized form, the discipline of literary criticism has a long and distinguished tradition (Habib, 2005). But ‘ordinary’ readers too have always talked about their reading experiences (Long, 2003; Rehberg Sedo, 2003). The written output of these reflections and discussions has been an important source for reading and reception studies. Proof of this importance is e.g. the existence of the Reading Experience Database (RED) that collects experiences of reading as documented in letters, memoirs and other historic material (Crone et al., 2011). Halsey (2009) e.g. shows how this database can help study changes in stylistic preferences over time.

One reason for the importance of written book responses is that they provide documentation of how works affect their readers: they show what elements of the reading experience readers consider important enough to write down and share with friends and fellow-readers. To some extent at least this will be determined by the elements of the book that were most significant to the reader and that he or she is most likely to remember. Unlike in earlier historic periods, this sort of evidence today is plentiful and researchers should take advantage of this. Spontaneous written responses to reading are not the only way of assessing the effects of (literary) reading. Experimental research (Miall, 2006) and other approaches have an important place. Today’s online book responses, however, are unique in that they are produced spontaneously by ordinary readers and have an ecological validity that other research data lack. (Which does, of course, not imply we should take everything that people write online at face value).

A second reason for the importance of written book responses is that their role as (co-)determiners, or at least predictors, of literary success is well-documented. In the wake of a large body of research on movie reviews (e.g. Liu, 2006), this was established for reviews on booksellers’ sites by (Chevalier and Mayzlin, 2006). For traditional (newspaper) reviews, their effects on long-term

success (canonization) have been shown in e.g. (Ekelund and Börjesson, 2002; Rosengren, 1987).

If reading responses are that important for the study of literature and its effects, it follows we need to understand them better. We need tools that can analyze their style, rhetorical structure, topics, and sentiment, and these tools should be sensitive to the many different sorts of readers, responses and response sites that form part of the landscape of online book discussion. We also need tools that can help us see relationships between the responses and the works that they respond to, in terms of topics and narrative (what characters and plot developments do reviewers respond to), as well as at higher (cognitive, emotional and moral) levels. An important step towards such tools is the creation of a representative corpus that can provide a test bed for tool development.

3 Related research

Online book discussion is a wide field that can be studied from many different angles. I discuss first a number of studies that do not use computational methods. Online book reviewing has often been discussed negatively in its relation to traditional reviews (McDonald, 2007; Pool, 2007). Certainly problematic aspects of online reviews are the possibilities of plagiarism and fraud (David and Pinch, 2006). Verboord (2010) uses a questionnaire to investigate the perceived legitimacy of internet critics. Online critics’ role in canonization was investigated in (Grafton, 2010). That online reviews do have an influence on books sales was established by (Chevalier and Mayzlin, 2006), and specifically for books by women and popular fiction in (Verboord, 2011). Many librarians have looked at what online book discussion sites can mean for the position of the library, library cataloguing and book recommendations (Pera and Ng, 2011; Pirmann, 2012). Online book discussion as an extension of the reading group is discussed in e.g. (Fister, 2005). A look at the whole field, from a genre perspective, is given in (Boot, 2011). Steiner (2010) looks specifically at Swedish weblogs; (Steiner, 2008) discusses Amazon reviews, as does (Domsch, 2009). Gutjahr (2002) sent out a survey to posters of Amazon reviews. Finally, (Miller, 2011) investigates how book blogs can

help develop the habits of mind required for literary reading.

Researchers that have used more or less sophisticated linguistic technology to investigate online book responses have done so with a number of different questions in mind. (Boot et al., 2012) sought to characterize responses from different site types based on word usage. Much effort has gone into the analysis of review sentiment, which has clear practical applications in marketing. (Taboada et al., 2011) use a lexicon-based approach; (Okanojima and Tsujii, 2005) a machine learning approach. (De Smedt and Daelemans, 2012a) create a Dutch sentiment lexicon based on reviews at an online bookseller. The helpfulness of online reviews has been investigated by e.g. (Tsur and Rappoport, 2009) while (Mukherjee and Liu, 2012) have modeled review comments. From an information retrieval perspective, the INEX social book search competition has explored the use of online reviews from Amazon and LibraryThing to create book recommendations (Koolen et al., 2012). A proposal for using text mining and discourse analysis techniques on pre-internet reviews is (Taboada et al., 2006). (Finn, 2011) used named entity recognition in reviews of a single writer in order to explore the ‘ideational network’ associated with her work.

It does not seem unfair to say that most of the computer-based linguistic research done into online book responses has been motivated by practical, if not commercial aims. Much of it was published in marketing journals. Computational linguistic research as a tool for understanding the variety of online book response is still at a very early stage of development.

4 A corpus of book responses

A corpus of book responses should present researchers with a varied, representative, and sufficiently large collection of book responses. It should not be a closed corpus but continue to grow. It should contain not just response texts but also include the metadata that describes and contextualizes the responses.

Varied: the responses should be taken from as wide a selection of sites as is possible. Sites are very different with regards to the active reviewers, their audience, the books that are discussed, the responses’ function and the explicit and im-

PLICIT expectations about what constitutes a proper response (Boot, 2011). Pragmatic aspects of the response (e.g. a response given on a weblog where the responder is the main author vs. a response in a forum where the responder is just one participant in a group discussion) obviously help determine both content and style of the response and tools that analyze responses should take account of these differences in setting.

Another respect in which variety is important is book genre. Much has been written about differences in book appreciation between e.g. readers of popular fiction and ‘high’ literature (Von Heydebrand and Winko, 1996). A response corpus should present researchers with a large body of responses from readers of a wide selection of genres (popular fiction, literature, non-fiction, essays, poetry, etc.), irrespective of its medium of publication (paper, e-book, online).

Representative: there is no need for this corpus to be strictly proportional with respect to site type or book genre. Still, it is important for all types and genres to be represented. Given the need to request permission from copyright holders, it will probably be impossible to achieve a truly representative corpus.

Sufficiently large: the required size of the corpus will depend on the sort of analysis that one tries to do. It is clear that analysis that goes beyond the collection level, e.g. at the book genre level, or at the level of individual reviewers, will need substantial amounts of text. A rule of thumb might be that collections should preferably contain more than a thousand responses and more than a million words.

Open: As new forms of computer-mediated communication continue to evolve, the ways of responding to and talking about books will also change. The corpus should facilitate research into these changes, and be regularly updated with collections from new site types.

Metadata: book response text acquires a large part of its meaning from its context. To facilitate research into many aspects of these responses it is important for the corpus to store information about that context. That information should include at least the site that the response was taken from, the response date, whatever can be known about the author of the response, and, if available, the book that the response responds to. Figure 1 shows the relevant entities.

We will not discuss the data model in detail. Sites can contain multiple collections of responses, with different properties. Some sites for instance contain both commissioned reviews and user reviews. Weblogs contains posts by the blog owner and responses to those posts. Book theme sites often carry review sections and discussion forums. When analyzing a response, it is important to be aware what section the response belongs to. Book responses can also be written in response to other posts, be it in a discussion forum, on Twitter, or on a book-based social networking site. Book responses can be tagged, and the tags may carry valuable information about book topics, book appreciation or other book information. Responses are written by persons, sometimes unknown, who may own a site (as with blogs) or be among many people active on a site, or perhaps on multiple sites. Reviewers sometimes write profile texts about themselves that also discuss their book preferences. On some sites (book SNS's, Twitter) reviewers may strike up friendships or similar relationships. Some sites also allow people to list the books they own and/or their favorite books. Finally, meaningful use of book level data will often require being able to group multiple versions (manifestations) of the same work.

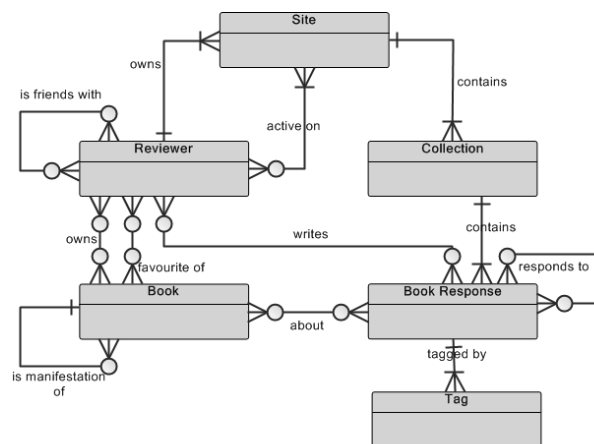


Figure 1. Book response corpus data model

For most collections, extracting the information carried by the respective entities mentioned is not a trivial task. Book shop review pages will probably contain an ISBN somewhere near the review, but forums probably will not and a tweet with an ISBN number is certainly unusual. And even if a response is ostensibly about book A, it may very

well also discuss book B. Reviewer information will also be hard to obtain, as many reviews (e.g. on booksellers' sites) are unsigned.

5 Pilot study

For a pilot study that explores the research potential of online book response, I have been collecting Dutch-language book responses from a number of sites. The size of the pilot corpus and its subcollections is given in table 1. The pilot corpus contains responses from a number of weblogs, from online review magazine 8Weekly, book-based social network site watleesjij.nu ('whatareyoureading.now'), book publicity, reviews and user reviews from thriller site Crimezone, a collection of print reviews (from multiple papers and magazines) about Dutch novelist Arnon Grunberg, print reviews from Dutch newspaper NRC and publicity from the NRC web shop. The collection should be extended with responses from other site types (e.g. forums, twitter, bookseller reviews) other book genres (e.g. fantasy, romance, poetry) and perhaps other text genres (e.g. book news, interviews).

Collection	Article genre	Response count	Word count (*1000)
8weekly	review	2273	1512
weblogs	blog post	6952	3578
watleesjij.nu	user review	28037	2515
crimezone book desc	publicity	3698	462
crimezone review	review	3696	1622
crimezone userrev	user review	9163	1537
grunberg	print review	196	187
NRC web shop	publicity	1345	198
NRC reviews	print review	1226	1133
Total		56586	12744

Table 1. Present composition of pilot corpus of responses

I have done a number of experiments in order to explore the potential for computational linguistic analysis of book responses.

5.1 Measure response style and approach using LIWC

As a first test, I investigated word usage in the book responses using LIWC (Pennebaker et al., 2007; Zijlstra et al., 2004). Figure 2 shows the usage of first person pronouns on the respective site types. The pattern conforms to what one would expect: on the book SNS *watleesjij.nu*, where readers give personal opinions, ‘I’ predominates, as it does in the Crimezone user reviews, and to a lesser extent in the weblogs. In the commissioned reviews both in print (NRC newspaper and Grunberg collection) and online (8Weekly) ‘we’ prevails, as reviewers have to maintain an objective stance. Interestingly, the Crimezone book descriptions manage to avoid first person pronouns almost completely.

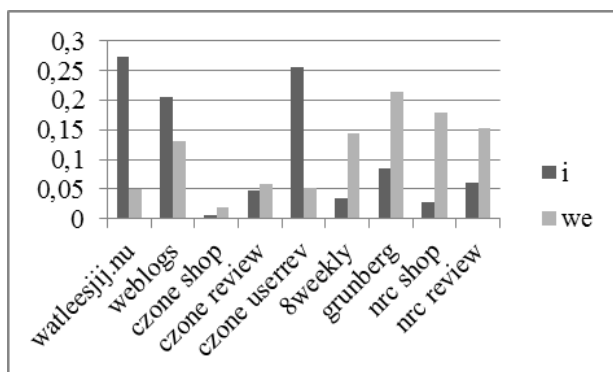


Figure 2. Normalized frequencies first person singular and first person plural pronouns

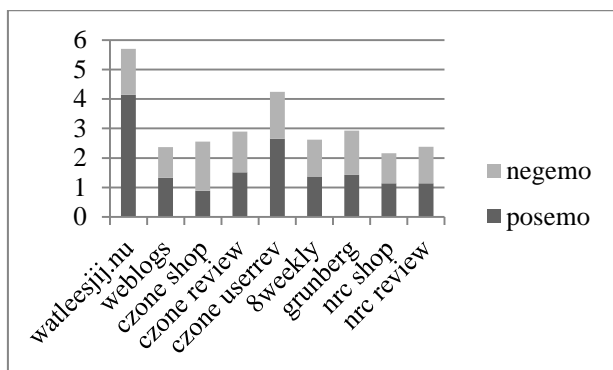


Figure 3. Positive and negative emotion word frequencies

A similar result appears when we chart positive and negative emotion words (Figure 3). Especially positive emotions are often expressed on *watleesjij.nu* and in the Crimezone user reviews. In this case the group of informal sites does not

include the weblogs, perhaps because the weblogs included in the pilot corpus are blogs at the intellectual end of the spectrum. Also interesting is the high proportion of negative emotion in the Crimezone book descriptions, perhaps because in the case of thrillers emotions like fear and anxiety can function as recommendations.

From these examples it is clear that word usage on the respective sites shows meaningful variation that will profit from further research. Investigation into these patterns at the level of individual reviewers (e.g. bloggers) should begin to show individual styles of responding to literature.

5.2 Site stylistic similarities

As a second test, I looked into writing style, asking whether the styles on the respective sites are sufficiently recognizable to allow meaningful clustering. For each of the collections, except for the weblogs, I created five files of 20000 words each and used the tools for computational stylometry described in (Eder and Rybicki, 2011) to derive a clustering, based on the 300 most frequent words. Figure 4 shows the results.

It is interesting to note that all except the *watleesjij.nu* (book SNS) samples are stylistically consistent enough to be clustered by themselves. It is even more interesting to note that the book descriptions from the NRC (newspaper) shop cluster with the descriptions taken from the Crimezone site, that the reviews in online magazine 8Weekly cluster with the printed reviews, and that the Crimezone reviews, commissioned and user-contributed, cluster with the *watleesjij.nu* reviews. This may be related to the fact that there are a large number of thriller aficionados on *watleesjij.nu*, or to Crimezone reviews being significantly different from traditional reviews. Again, this seems a fruitful area for further investigation, only possible in the context of a large corpus containing different text types.

In order to exclude the possibility that this clustering is based on content words (e.g. words related to crime), I repeated the experiment using bi-grams of the words’ POS-tags, as derived by the Pattern toolset (De Smedt and Daelemans, 2012b). The resulting figure, not reproduced here, is very similar to Figure 4. This result leads to another question: what sort of syntactic construc-

tions are specific to which site types? And can we connect these stylistic differences to the approach to literature that these sites take?

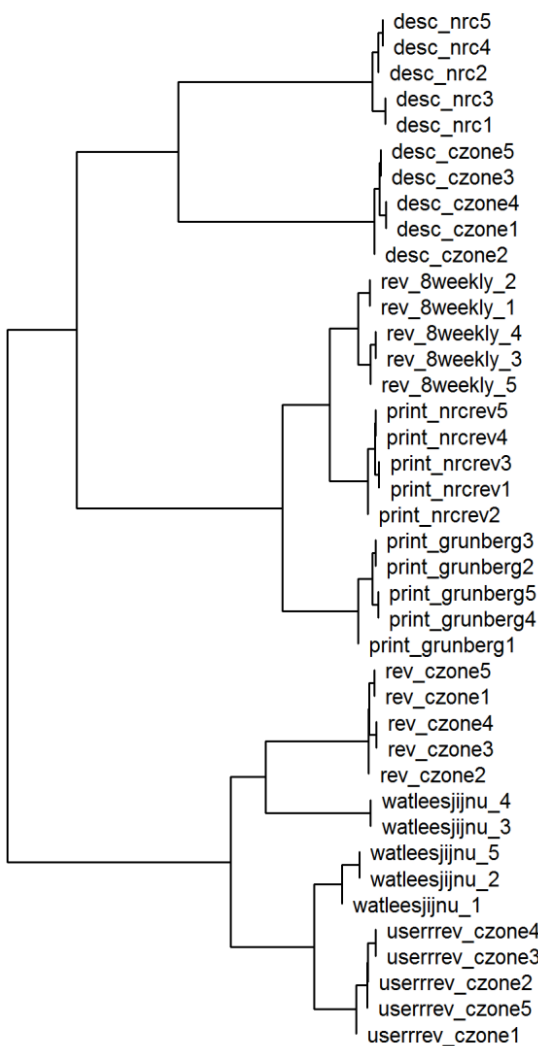


Figure 4. Clustering of 20000-word review texts based on 300 most frequent words.

5.3 Response sentiment analysis

In a third experiment, I applied the sentiment lexicon embedded in the Pattern toolset to the responses in those collections that include ratings. I predict a positive rating (i.e. above or equal to the collection median) when the sentiment as measured by Pattern is above 0.1, and compute precision, recall and F1-value for this prediction (see Figure 5). Results on the book SNS watleesjij.nu are similar to the results reported by (De Smedt and Daelemans, 2012a) for reviews from

bookseller bol.com, perhaps because the responses on the two sites are similar. As expected, the results are considerably worse for the longer reviews on 8Weekly and NRC. That precision should be as high as .84 for the Crimezone reviews is somewhat of a mystery.

While it is not unexpected that the sentiment prediction quality should be higher for the sites with simpler reviews, this does imply a challenge for researchers of sentiment analysis. Without accurately gauging response sentiment (and many other response properties) measuring literary impact from responses will remain illusory.

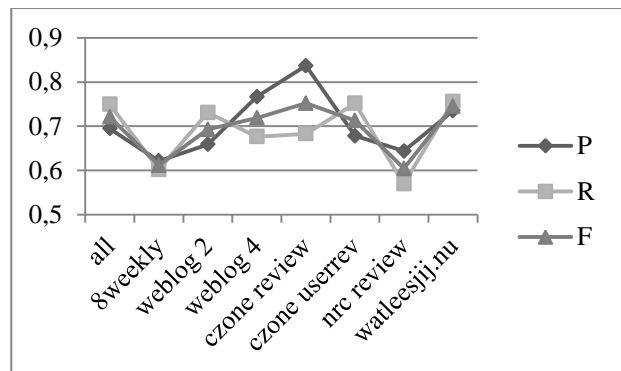


Figure 5. Prediction of positive or negative rating: precision, recall and F-score

5.4 Semantic similarities between book fragments and responses

A final experiment is based on the assumption that the semantics of book response texts to some extent reflect the semantics of the books they respond to. If that is true, it should be possible to determine the chapters that most impressed readers by comparing the book's and the reviews' semantic content. In order to test the assumption, I used Latent Semantic Analysis (LSA) (Landauer et al., 2007; Řehůřek and Sojka, 2010) to measure the distances between 400-word fragments taken from the novel *Tirza* by Dutch novelist Arnon Grunberg and 60 reviews of the book taken from book SNS watleesjij.nu. In order to compensate for potential similarities between book fragments and any reviews, rather than with reviews specifically of this book, I also measured semantic distances between the book's fragments and a set of random reviews from the same site, and subtracted those from the distances with the *Tirza* reviews. In order to test how these distances relate

to the book's content, I computed LIWC scores for the fragments and then correlations between these LIWC scores and the LSA distances. For e.g. LIWC category 'family', a very important subject for this book, the correlation is positive and highly significant (.34, $p < .0001$).

Further experimentation with other books, other review collections and other LSA models is clearly needed. It is too early to say whether LSA indeed offers a viable approach for determining the book fragments most closely related to review texts, but this is clearly a promising result. Being able to connect measurable aspects of books with impact in reviews would help us understand how books affect their readers.

6 Conclusion

This paper adopts a broad conception of the object of literary studies, taking it to include the individual and social responses that literature elicits. I argued here that the (plentifully available) online book responses are important to literary studies, both as evidence (because they document the reception of literary works) and as objects (because they help determine works' short and long term popularity). If only because of the numbers of these responses, we need computational linguistic tools in order to analyze and understand them. Because the responses published on the various response platforms are in many respects very different, potential tools would need to be developed with these differences in mind. A good way to ensure this is to create an appropriately large and representative corpus of online book response. On the basis of a Dutch pilot corpus we saw that existing linguistic tools can reveal some of the differences between the respective platforms. They are currently unable, however, to perform any deeper analysis of these differences, let alone a deeper analysis of the relations between responses and books.

Naturally, written book response can only inform us about the reading experience of those that take the trouble of writing down and publishing their response. Even though those who provide book response are by no means a homogeneous group, it is clear that the proposed corpus would necessarily be selective, and should not be our only method of studying reader response. This is less of an issue when studying how books become

popular and eventually canonized, as those who don't participate in the discussions will, for that very reason, be less influential.

With these caveats, there are a number of areas that a corpus of online book response would help investigate. Among these are:

- the responses themselves and their respective platforms: what language is used, what topics are discussed, what is their structure? What do they reveal about the literary norms that (groups of) readers apply?
- the relations between responses: we should be able to answer the questions about influence. What sort of discussions are going on about literature on which platforms? Which participants are most influential? Can response styles reveal these influences?
- what the responses show about the reading experience: we'd like to know how books (both books in general and specific books) affect people, what attracts people in books, what they remember from books, what they like about them, etc. What passages do they quote from the books they respond to? What characteristic words do they adopt?
- what the responses show about readers: as the corpus should facilitate selection by responder, we should be able to investigate the role of the reader in book response. Do responders' writing styles predict their ratings? Do people who like, say, James Joyce dislike science fiction? And can their book responses tell us why?

Many of these phenomena are interesting at multiple levels. They are interesting at the level of the individual reader, for whom reading in general and specific books are important. They are interesting at a sociological level, as discussions help determine books' popularity or even canonization. Finally, at the level of the book, study of book responses can show what readers, individually and in groups, take away from a book. In this respect especially, study of book responses is a necessary complement to study of the literary text.

References

- Boot, Peter. 2011. Towards a Genre Analysis of Online Book Discussion: socializing, participation and publication in the Dutch booksphere. *Selected Papers of Internet Research* IR 12.0.
- Boot, Peter, Van Erp, Marieke, Aroyo, Lora, and Schreiber, Guus. 2012. The changing face of the book review. Paper presented at *Web Science 2012*, Evanston (IL).
- Boyd, Brian. 2009. *On the origin of stories: Evolution, cognition, and fiction*. Cambridge MA: Harvard University Press.
- Chevalier, Judith A., and Mayzlin, Dina. 2006. The effect of word of mouth on sales: Online book reviews. *Journal of Marketing Research* 43:345-354.
- Crone, Rosalind, Halsey, Katry, Hammond, Mary, and Towheed, Shafquat. 2011. The Reading Experience Database 1450-1945 (RED). In *The history of reading. A reader*, eds. Shafquat Towheed, Rosalind Crone and Katry Halsey, 427-436. Oxon: Routledge.
- David, Shay, and Pinch, Trevor. 2006. Six degrees of reputation: The use and abuse of online review and recommendation systems. *First Monday* 11.
- De Smedt, Tom, and Daelemans, Walter. 2012a. "Vreselijk mooi!" (terribly beautiful): A Subjectivity Lexicon for Dutch Adjectives. Paper presented at *Proceedings of the 8th Language Resources and Evaluation Conference (LREC'12)*.
- De Smedt, Tom, and Daelemans, Walter. 2012b. Pattern for Python. *The Journal of Machine Learning Research* 13:2031-2035.
- Domsch, Sebastian. 2009. Critical genres. Generic changes of literary criticism in computer-mediated communication. In *Genres in the Internet: issues in the theory of genre*, eds. Janet Giltrow and Dieter Stein, 221-238. Amsterdam: John Benjamins Publishing Company.
- Eder, Maciej, and Rybicki, Jan. 2011. Stylometry with R. In *Digital Humanities 2011: Conference Abstracts*, 308-311. Stanford University, Stanford, CA.
- Ekelund, B. G., and Börjesson, M. 2002. The shape of the literary career: An analysis of publishing trajectories. *Poetics* 30:341-364.
- Finn, Edward F. 2011. *The Social Lives of Books: Literary Networks in Contemporary American Fiction*, Stanford University: PhD.
- Fister, Barbara. 2005. Reading as a contact sport. *Reference & User Services Quarterly* 44:303-309.
- Goldstein, Philip, and Machor, James L. 2008. *New directions in American reception study*. New York: Oxford University Press, USA.
- Grafton, Kathryn. 2010. *Paying attention to public readers of Canadian literature: popular genre systems, publics, and canons*, University of British Columbia: PhD.
- Gutjahr, Paul C. 2002. No Longer Left Behind: Amazon.com, Reader-Response, and the Changing Fortunes of the Christian Novel in America. *Book History* 5:209-236.
- Habib, M. A. R. 2005. *A history of literary criticism: from Plato to the present*. Malden, MA: Blackwell.
- Halsey, Katie. 2009. 'Folk stylistics' and the history of reading: a discussion of method. *Language and Literature* 18:231-246.
- Koolen, Marijn, Kamps, Jaap, and Kazai, Gabriella. 2012. Social Book Search: Comparing Topical Relevance Judgements and Book Suggestions for Evaluation. In *CIKM'12, October 29–November 2, 2012*. Maui, HI, USA.
- Landauer, T. K., McNamara, D. S., Dennis, S., and Kintsch, W. 2007. *Handbook of latent semantic analysis*: Lawrence Erlbaum.
- Liu, Yong. 2006. Word of Mouth for Movies: Its Dynamics and Impact on Box Office Revenue. *Journal of Marketing* 70:74-89.
- Long, Elizabeth. 2003. *Book clubs: Women and the uses of reading in everyday life*. Chicago: University of Chicago Press.
- McDonald, Rónán. 2007. *The death of the critic*. London, New York: Continuum International Publishing Group.
- Miall, David S. 2006. *Literary reading: empirical & theoretical studies*. New York: Peter Lang Publishing.
- Miller, Donna L. 2011. *Talking with Our Fingertips: An Analysis for Habits of Mind in Blogs about Young Adult Books*, Arizona State University: PhD.
- Mukherjee, Arjun, and Liu, Bing. 2012. Modeling Review Comments. In *Proceedings of 50th Annual Meeting of Association for Computational Linguistics (ACL-2012)*. Jeju (Korea).
- Okanohara, Daisuke, and Tsujii, Jun'ichi. 2005. Assigning polarity scores to reviews using machine learning techniques. *Natural Language Processing–IJCNLP 2005*:314-325.
- Pennebaker, J. W., Booth, R. J., and Francis, M. E. 2007. *Linguistic Inquiry and Word Count (LIWC2007)*. Austin, TX.
- Pera, Maria Soledad, and Ng, Yiu-Kai. 2011. *With a Little Help from My Friends: Generating*

- Personalized Book Recommendations Using Data Extracted from a Social Website. Paper presented at *Web Intelligence and Intelligent Agent Technology (WI-IAT)*, 2011.
- Pirmann, Carrie. 2012. Tags in the Catalogue: Insights From a Usability Study of LibraryThing for Libraries. *Library Trends* 61:234-247.
- Pool, Gail. 2007. *Faint praise: the plight of book reviewing in America*. Columbia, MO: University of Missouri Press.
- Rehberg Sedo, DeNel. 2003. Readers in Reading Groups. An Online Survey of Face-to-Face and Virtual Book Clubs. *Convergence* 9:66-90.
- Řehůřek, Radim, and Sojka, Petr. 2010. Software framework for topic modelling with large corpora. Paper presented at *Proceedings of LREC 2010 workshop New Challenges for NLP Frameworks*, Valletta, Malta.
- Rosengren, Karl Erik. 1987. Literary criticism: Future invented. *Poetics* 16:295-325.
- Steiner, Ann. 2008. Private Criticism in the Public Space: Personal writing on literature in readers' reviews on Amazon. *Participations* 5.
- Steiner, Ann. 2010. Personal Readings and Public Texts: Book Blogs and Online Writing about Literature. *Culture unbound* 2:471-494.
- Taboada, Maite, Gillies, Mary Ann, and McFetridge, Paul. 2006. Sentiment Classification Techniques for Tracking Literary Reputation. In *Proceedings of LREC 2006 Workshop "Towards Computational Models of Literary Analysis"*.
- Taboada, Maite, Brooke, Julian, Tofiloski, Milan, Voll, Kimberly, and Stede, Manfred. 2011. Lexicon-based methods for sentiment analysis. *Computational Linguistics* 37:267-307.
- Tsur, Oren, and Rappoport, Ari. 2009. Revrank: A fully unsupervised algorithm for selecting the most helpful book reviews. Paper presented at *International AAAI Conference on Weblogs and Social Media*.
- Van Peer, Willie. 2008. Introduction. In *The quality of literature: linguistic studies in literary evaluation*, 1-14. Amsterdam: John Benjamins Publishing Co.
- Verboord, Marc. 2010. The Legitimacy of Book Critics in the Age of the Internet and Omnivorousness: Expert Critics, Internet Critics and Peer Critics in Flanders and the Netherlands. *European Sociological Review* 26:623-637.
- Verboord, Marc. 2011. Cultural products go online: Comparing the internet and print media on distributions of gender, genre and commercial success. *Communications* 36:441-462.
- Von Heydebrand, Renate, and Winko, Simone. 1996. *Einführung in die Wertung von Literatur: Systematik, Geschichte, Legitimation*. Paderborn: Schöningh.
- Zijlstra, Hanna, van Meerveld, Tanja, van Middendorp, Henriët, Pennebaker, James W., and Geenen, Rinie. 2004. De Nederlandse versie van de 'Linguistic Inquiry and Word Count' (LIWC). *Gedrag & Gezondheid* 32:271-281.