# The 'Tiroler Soldaten-Zeitung' and its Authors. A Computer-Aided Search for Robert Musil

**Salgaro, Massimo**

massimo.salgaro@univr.it
University of Verona, Italy

**Rebora, Simone**

simone.rebora@univr.it
University of Verona, Italy

**Lauer, Gerhard**

gerhard.lauer@unibas.ch
University of Basel, Switzerland

**Herrmann, J. Berenike**

berenike.herrmann@unibas.ch
University of Basel, Switzerland

Robert Musil, one of the most important authors of the twentieth-century German-written literature, fought in the Austrian army at the Italian front. During the First World War, between 1916 and 1917, Musil was chief editor of the *Tiroler Soldaten-Zeitung* in Bozen. This activity has always been a philological problem for Musil scholars, who have not been able to attribute with certainty a range of texts to the author. However, their identification is fundamental in the study of his political thinking. With this paper, we present a new approach, that combines historical and philological research with stylometric methods.

The starting point for the determination of possible authorship is the screening of previous attempts. The number of articles attributed to Musil has so far varied extensively:

| Attribution proposed by | Number of TSZ articles attributed to Musil |
|---|---|
| (Dinklage 1960) | 3 |
| (Roth 1972) | 19 |
| (Corino 1973, 2003, and 2010) | 8 |
| (Arntzen 1980) | 22 |
| (Fontanari / Libardi 1987) | 36 |
| (Amann *et al.* 2009) | 36 |

We have limited our test set to the 38 TSZ articles listed by (Schaunig 2014), for which Musil's authorship has been proposed at least once. The major problem for carrying out a stylometric analysis on this corpus is text length. As demonstrated by recent research, the minimum length for a reliable authorship attribution is around 5,000 words (see Eder 2015). However, the average length of the 38 disputed TSZ articles is slightly below 1,000 words (see Figure 1). As a possible solution for this issue, we decided to develop a combinatory design that analyzes longer chunks composed by the juxtaposition of single texts. To reduce the number of combinations, we excluded the nine shortest texts (below 500 word), together with the only text attributed to Musil on solid philological ground (see Corino 1973). This leaves us with a corpus of 28 texts, already digitized by (Amann *et al.* 2009). The optimal configuration was obtained by combining groups of 6 texts. This permutation generated 376,740 text chunks with an average length of N=6,963 words and a standard deviation of 909 words.
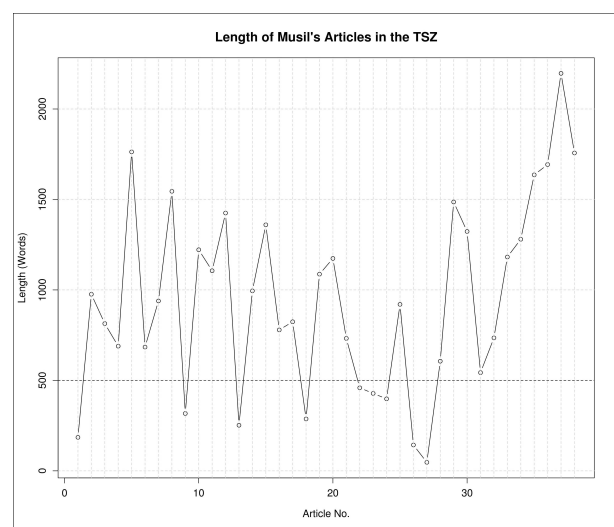


*Figure 1. Test set composition*

As for the composition of the training set, we drew both on the "impostors method" (see Koppel / Winter 2014) and on historiographical research. Following (Juola 2015), we fixed the number of "impostors" to a minimum of three: Franz Blei, Franz Kafka, and Stefan Zweig. Subsequently, we selected three authors suggested by (Urbaner 2001) as possible TSZ collaborators: Marie delle Grazie, Hugo Salus, and Albert Ritter (his texts were not available in digitized format, so we OCRed and manually refined them). The training set was then completed by a selection of articles published by Musil in various journals between 1911 and 1919. For each author, the retrieved material was subdivided in three text chunks with a length comprised between 6,000 and 8,000 words: the training set was thus composed by 21 text chunks (see Figure 2).
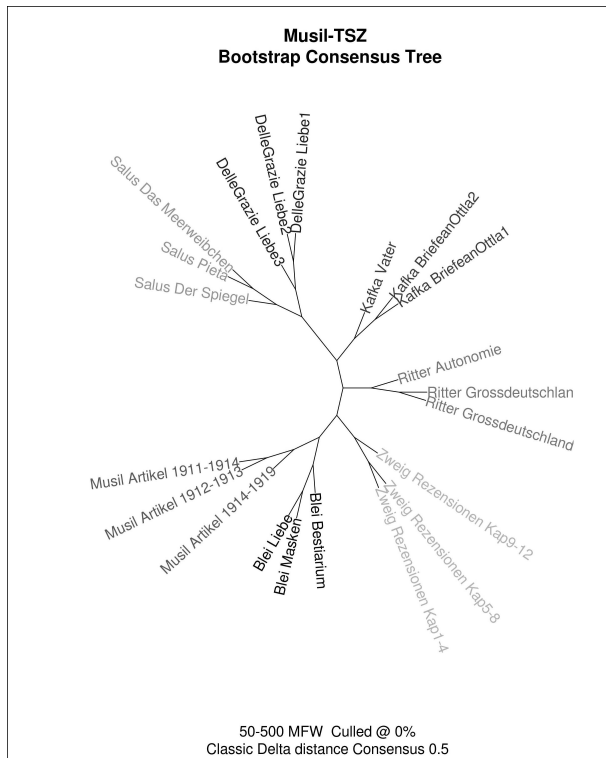
*Figure 2. Training set*



*Figure 3. Combinatory design results*

The analysis was carried out using the R package *Stylo* (see Eder / Rybicki / Kestemont 2016). For each iteration, the distances between test set and training set were saved in the tabular form provided by the package. At the end of the process, mean values were calculated by sub-grouping the combinations by each TSZ text. Notwithstanding the employment of a high-standard computational power (provided by GWDG, University of Göttingen), a first experiment using 50–500 most frequent words (MFW) and Burrows's Delta distance took more than one week to be completed (see Figure 3). However, when repeating the experiment with only one-tenth of the combinations (i.e. 37,674 iterations, randomly selected), results were rather identical (see Figure 4) and the process took less than one day. When the experiment was repeated without any combination, results were extremely noisier (see Figure 5), thus confirming that the combinatory design was able to better discern authorial signals.
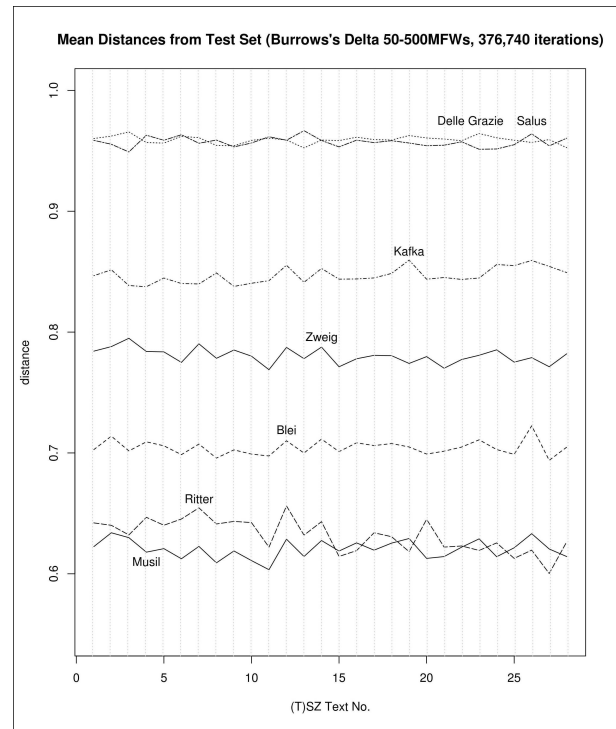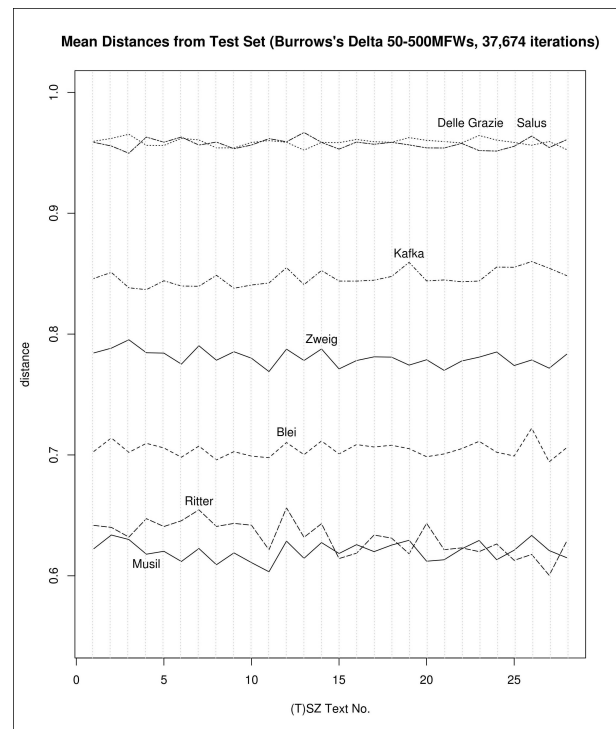


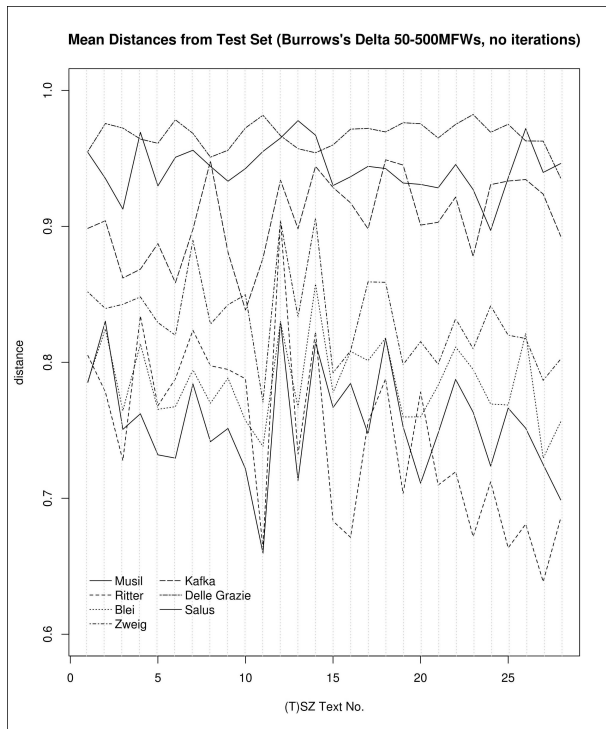*Figure 4. Combinatory design results (one-tenth iterations)*

*Figure 5. Results without combinatory design*

To validate the results, the experiment has been repeated with 16 different configurations, by combining Eder's Delta, Burrow's Delta, Canberra, and Cosine distances with 10–100, 20–200, 50–500, and 100–1,000 MFW. In all configurations, Ritter and Musil are the only authors disputing the authorship of the TSZ articles. This evidence has been corroborated by the discovery of a document in the *Kriegsarchiv* in Wien, which confirms that Albert Ritter was part of the TSZ editorial team (see Figure 6).
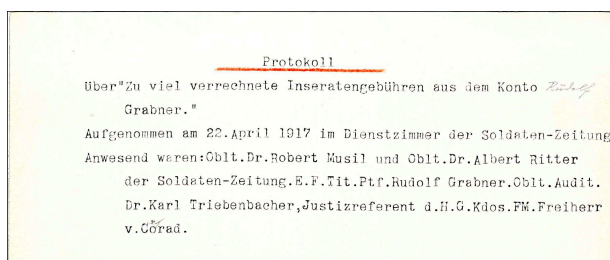


*Figure 6. Source: Kriegsarchiv, Wien*

Final results have been synthetized here:

| TSZ articles' titles and dates of publication | Agreement between classifiers on Musil's authorship |
|---|---|
| 1. „Kameraden arbeitet mit!" (6. 8. 1916) | 100,00% |
| 2. „Bin ich ein Österreicher?" (20. 8. 1916) | 87,50% |
| 3. „Herr Tüchtig und Herr Wichtig" (27. 8. 1916) | 81,25% |
| 4. „Das Schlagwort" (27. 8. 1916) | 100,00% |
| 5. „Die Erziehung zum Staat" (3. 9. 1916) | 100,00% |
| 6. „Bauernleben" (1. 10. 1916) | 100,00% |
| 7. „Sonderbare Patrioten" (15. 10. 1916) | 100,00% |
| 8. „Noch einmal Bauernleben" (29. 10. 1916) | 100,00% |
| 9. „Opportunität" (12. 11. 1916) | 100,00% |
| 10. „Eine gute persönliche Beziehung" (26. 11. 1916) | 100,00% |
| 11. „Eine österreichische Kultur" (10. 12. 1916) | 100,00% |
| 12. „Der Nörgler und der neue Österreicher" (17. 12. 1916) | 100,00% |
| 13. „Das Kompromiß" (24. 12. 1916) | 100,00% |
| 14. „Heilige Zeit" (31. 12. 1916) | 100,00% |
| 15. „Zentralismus und Föderalismus" (7. 1. 1917) | 68,75% |
| 16. „Föderalismus oder Zentralismus" (14. 1. 1917) | 68,75% |
| 17. „Zu Milde und zu Wilde" (11. 2. 1917) | 93,75% |
| 18. „Neu-Altösterreichisches" (25. 2. 1917) | 87,50% |
| 19. „Ist die »österreichische Frage« schwierig?" (4. 3. 1917) | 62,50% |
| 20. „Seiner Hochwohlgeboren!" (4. 3. 1917) | 100,00% |
| 21. „Luxussteuern" (4. 3. 1917) | 93,75% |
| 22. „Positive Ziele" (11. 3. 1917) | 81,25% |
| 23. „Der Frieden versprochen!" (18. 3. 1917) | 68,75% |
| 24. „Das Staatsprogramm der Deutschen" (18. 3. 1917) | 87,50% |
| 25. „Wehe dem Staatsmann!" (25. 3. 1917) | 68,75% |
| 26. „Der Frieden und die Zukunft" (1. 4. 1917) | 62,50% |

| 27. „Presse und Krieg" (8. 4. 1917) | 68,75% |
|---|---|
| 28. „Vermächtnis" (15. 4. 1917) | 100,00% |

A general trend is evident: while, for the articles published in 1916, Musil's authorship is almost unquestionable, many more doubts emerge with the articles published in 1917. In no case, however, Ritter's signal becomes dominant. Notwithstanding the high margins of uncertainty, these results are to be considered as significant for multiple reasons. First, the combinatory design, while having shown the dominance of Musil's signal throughout the test set, may have overshadowed different, minor signals. Second, it should be considered the fact that Musil, in the role of chief editor, may have altered many articles in the journal, thus intermixing his authorial signal with those of others. All this considered, further research is advisable, while the focus should be shifted towards the texts on which classifiers disagree.

Among possible future developments of the research, is the definition of new training sets to validate the results and an expansion of the test set. Both these developments, however, will require an extensive digitization effort: most of the useful texts, in fact, are not available in a clean plain-text format. In addition, other software should be tested on the already defined corpus, from JGAAP (see Juola *et al.* 2008) to the CLEF/PAN software (see Stamatatos *et al.* 2014), focusing specifically on different methods for authorship attribution, from lower-level features such as character n-grams (see Halvani *et al.* 2016), to higher-level features such as syntactic labels (see Hirst / Feiguina 2007), taking into consideration also machine-learning techniques (see Jockers / Witten 2010). With our study, we hope to have cast the groundwork for a research that can have long-lasting consequences on the history of German literature, confirming at the same time how quantitative methods are not in opposition, but complementary to qualitative analysis (see Herrmann 2018).

# Bibliography

**Amann, Klaus / Corino, Karl / Fanta, Walter** (2009): *Robert Musil, Klagenfurter Ausgabe.* Klagenfurt: Robert Musil-Institut der Universität Klagenfurt.

**Arntzen, Helmut** (1980): *Musil-Kommentar sämtlicher zu Lebzeiten erschienener Schriften außer dem Roman "Der Mann ohne Eigenschaften".* München: Winkler.

**Corino, Karl** (1973): "Robert Musil, Aus der Geschichte eines Regiments", in: Studi Germanici 11: 109–115.

**Corino, Karl** (2003): *Robert Musil: eine Biographie*, Reinbek bei Hamburg: Rowohlt.

**Corino, Karl** (2010): "Klaviersonnen über Schluchten des Gemüts. Robert Musil und die Musik", in: *Das Plateau* 120: 4–21.

**Dinklage, Karl** (1960): *Robert Musil. Leben, Werk, Wirkung*, Zürich.

**Eder, Maciej / Kestemont, Mike / Rybicki, Jan** (2016): "Stylometry with R: a package for computational text analysis", in: *R Journal* 8(1): 107–121.

**Eder, Maciej** (2015): "Does size matter? Authorship attribution, small samples, big problem", in: *Digital Scholarship in the Humanities* 30(2): 167–182.

**Fontanari, Alessandro / Libardi, Massimo** (1987): *La guerra parallela*. Trento: Reverdito.

**Halvani, Oren / Winter, Christian / Pflug, Anika** (2016): "Authorship verification for different languages, genres and topics", in: *Digital Investigation* 16: 33–43.

**Herrmann, J. Berenike** (2018): "In test bed with Kafka. Introducing a mixed-method approach to digital stylistics", in: *Digital Humanities Quarterly* [in press].

**Hirst, Graeme / Feiguina, Ol'ga** (2007): "Bigrams of syntactic labels for authorship discrimination of short texts", in: *Literary and Linguistic Computing* 22(4): 405–417.

**Jockers, Mattew / Witten, Daniela** (2010): "A comparative study of machine learning methods for authorship attribution", in: *Literary and Linguistic Computing* 25(2): 215–223.

**Juola, Patrick / Noecker, John / Ryan, Mike / Zhao, Mengjia** (2008): "JGAAP3.0 – authorship attribution for the rest of us", in: *Digital Humanities 2008: Book of Abstracts*: 250–251.

**Juola, Patrick** (2015): "The Rowling Case: A Proposed Standard Analytic Protocol for Authorship Questions", in: *Digital Scholarship in the Humanities* 30: 100–113.

**Koppel, Moshe / Winter Yaron** (2014): "Determining if two documents are by the same author", in: JASIST 65(1): 178-187.

**Roth, Marie-Louise** (1972): *Robert Musil. Ethik und Ästhetik*. München: List.

**Schaunig, Regina** (2014): *Der Dichter im Dienst des Generals. Robert Musils Propagandaschriften im ersten Weltkrieg*. Klagenfurt: Kitab.

**Stamatatos, Efstathios / Daelemans, Walter / Verhoeven, Ben / Potthast, Martin / Stein, Benno / Juola, Patrick / Sanchez-Perez, Miguel A. / Barrón-Cedeño, Alberto** (2014): "Overview of the Author Identification Task at PAN 2014 ", in: *CLEF 2014 (Working Notes)*: 877-897.

**Urbaner, Roman** (2001): "'... daran zugrunde gegangen, daß sie Tagespolitik treiben wollte'? Die ,(Tiroler) Soldaten-Zeitung' 1915-1917", in: *eForum zeitGeschichte* 3/4. www.eforum-zeitgeschichte.at [accessed 14.09.2017]