

---

# HOW DOES BURROWS' DELTA WORK ON MEDIEVAL CHINESE POETIC TEXTS?

---

A PREPRINT

 **Boris Orekhov\***

School of Linguistics

HSE University,

Institute of Russian Literature (Pushkin House)

Russian Academy of Sciences

borekhov@hse.ru

July 12, 2024

## ABSTRACT

Burrows' Delta was introduced in 2002 and has proven to be an effective tool for author attribution. Despite the fact that these are different languages, they mostly belong to the same grammatical type and use the same graphic principle to convey speech in writing: a phonemic alphabet with word separation using spaces. The question I want to address in this article is how well this attribution method works with texts in a language with a different grammatical structure and a script based on different principles. There are fewer studies analyzing the effectiveness of the Delta method on Chinese texts than on texts in European languages. I believe that such a low level of attention to Delta from sinologists is due to the structure of the scientific field dedicated to medieval Chinese poetry. Clustering based on intertextual distances worked flawlessly. Delta produced results where clustering showed that the samples of one author were most similar to each other, and Delta never confused different poets. Despite the fact that I used an unconventional approach and applied the Delta method to a language poorly suited for it, the method demonstrated its effectiveness. Tang dynasty poets are correctly identified using Delta, and the empirical pattern observed for authors writing in European standard languages has been confirmed once again.

**Keywords** tang poetry · burrows's delta · author attribution · chinese language

## 1 Introduction

Burrows' Delta was introduced in 2002 [Burrows, 2002] and has proven to be an effective tool for author attribution. The psychological and linguistic foundations of this method are not entirely clear, but numerous successful tests in various languages have confirmed that it is a viable method for calculating intertextual distance, which correlates well with a specific author's relationship to a text. If several texts are written by the same person, the Delta distance between them is smaller than between texts written by different authors.

Although this rule does not work in 100% of cases [Skorinkin and Orekhov, 2023], it covers most cases that scholars in philology are interested in. The number of studies in which Burrows' Delta becomes a tool for solving attribution problems is growing. There are two reasons for this: first, the effectiveness of the method itself, repeatedly tested on texts in different languages. Second, the low entry threshold for researchers, provided by the well-crafted stylo package for the R language [Eder et al., 2016], which generally follows a zero-coding strategy.

The effectiveness of the method has been proved on the material of different languages, both new and ancient: English [Hoover, 2004], Old English [Garcia and Martin, 2006], German [Jannidis and Lauer, 2014], Spanish [Hernández-Lorenzo and Byszuk, 2023], Italian [Savoy, 2018, Rybicki, 2018], Polish [Rybicki and Heydel, 2013], Russian [Skorinkin and Bonch-Osmolovskaya, 2016], Arabic [AbdulRazzaq and Mustafa, 2014], and others. There is a paper comparing the quality of attribution in different languages: Latin, Polish, English, German [Eder, 2011].

Despite the fact that these are different languages, they mostly (with the exception of Arabic) belong to the same grammatical type and use the same graphic principle to convey speech in writing: a phonemic alphabet with word separation using spaces. The question I want to address in this article is how well this attribution method works with texts in a language with a different grammatical structure and a script based on different principles.

There are fewer studies analyzing the effectiveness of the Delta method on Chinese texts than on texts in European languages. It is worth noting [du2, 2016] that those studies focus on modern Chinese, whereas I am interested in medieval Chinese poetry. This is a language without inflection and a script system without spaces as word delimiters. With such a grammatical structure, the distribution of function words is different from what is typical in European-standard languages [Xiao, 2008, Liu et al., 2017]. Furthermore, in the writing systems familiar to Europeans, spaces separate words from each other. In Chinese, tokenization presents a complex problem [Huang and Wu, 2018], so it is important to check if Delta can be used by focusing only on the obvious text units, the characters, without resorting to complex algorithms that can introduce errors into text processing. Although the use of letter ngrams for Delta attribution is also found in European languages, typically sequences of several characters are used. Since in Chinese, the significance of a single character is higher than that of an individual letter in European languages, we will attempt to rely on individual characters rather than sequences.

I believe that such a low level of attention to Delta from sinologists is due to the structure of the scientific field dedicated to medieval Chinese poetry. Medieval Chinese literature is very different from medieval European literature, which had many gaps. The Chinese tradition of this period was highly regulated [Sturgeon, 2018], as was much of Chinese society, and it was built on the concept of authority, so texts were well documented. As a result, the material generally poses far fewer problems in terms of determining authorship. Authorship determination is the area where Delta traditionally shows the most significant results, which is why it is not in demand among sinologists specializing in medieval studies, as the authorship of texts is well known.

## 2 Method and Data

There is no need to recount the calculation algorithm for computing Delta once again. This has been done dozens of times in the most well-known and authoritative publications. The intertextual distance is calculated using the formula:

$$\Delta = \sum_{i=1}^n \frac{|z(x_i) - z(y_i)|}{n} \quad (1)$$

Thus, we determine the distance between each and every text in the research sample. Since we do not have a reliable method for tokenizing the text, we are forced to work with individual characters. The Stylo package has an option that allows switching from words to letter n-grams. This means we treat each character as a separate letter, which is an intermediate approach between alphabetic scripts and the script characteristic of the Chinese language, where one word can be represented by multiple characters.

### 2.1 Data

As the source of texts, the collection stored in the repository at [snowtraces, 2020] was used. It contains a digitized version of the "Complete Tang Poems" (or *Quan Tangshi*), which is the largest collection of Tang poetry. Thus, this study examines the effectiveness of the Delta method for determining authorship in Tang dynasty Chinese poetry.

I collected all texts of a single author together, resulting in 2537 poets. This is too many for a visual analysis of the results. Therefore, I further worked only with the 20 most prolific poets of the Tang era. Here is their list and the number of characters for each, including characters, spaces, line breaks, and punctuation marks:

1. 白居易 229346
2. 杜甫 129391
3. 李白 103294
4. 元稹 80968
5. 韓愈 60013
6. 劉禹錫 58224
7. 貫休 48836
8. 齊己 46381

9. 陸龜蒙 45104
10. 韋應物 40776
11. 孟郊 40582
12. 李商隱 39589
13. 皎然 38349
14. 劉長卿 35301
15. 皮日休 34357
16. 杜牧 33483
17. 王建 31607
18. 姚合 30887
19. 錢起 29874
20. 許渾 29775

For stylometry, volume is important, so I combined the poems of one poet into several large samples to compare them with each other using Delta. However, the results can depend on which specific poems we combine. For example, one combination of poems might yield one version of stylometric distribution, while another combination might yield a different version. Therefore, it is necessary to try different combinations of poems within the sample. I did the following: I randomly mixed the order of poems for one author, split the volume in half, and presented each half as separate texts (samples), repeating this five times to create five test corpora. The code for these manipulations is provided in this repository [Orekhov, 2024].

Next, I applied Delta to all these corpora using the Stylo package, configuring the package to work with characters instead of words and selecting the 100 most frequent characters. Here is the list of the 10 most frequent characters from the analysis:

1. 不
2. 人
3. 無
4. 一
5. 日
6. 山
7. 風
8. 有
9. 何
10. 來

### 3 Results

In all 5 versions of the corpus, clustering based on intertextual distances worked flawlessly. Delta produced results where clustering showed that the samples of one author were most similar to each other, and Delta never confused different poets. We present only one figure 1 as an illustration, as the others repeat the same result.

Let's provide a table 1 with a portion of the distances. The full table can be viewed in the data publication [Orekhov, 2024].

The heatmap (Fig. 2) of distances also records the smallest distances between samples of the same poet.

### 4 Conclusion

Despite the fact that I used an unconventional approach and applied the Delta method to a language poorly suited for it, the method demonstrated its effectiveness. Tang dynasty poets are correctly identified using Delta, and the empirical pattern observed for authors writing in European standard languages has been confirmed once again.

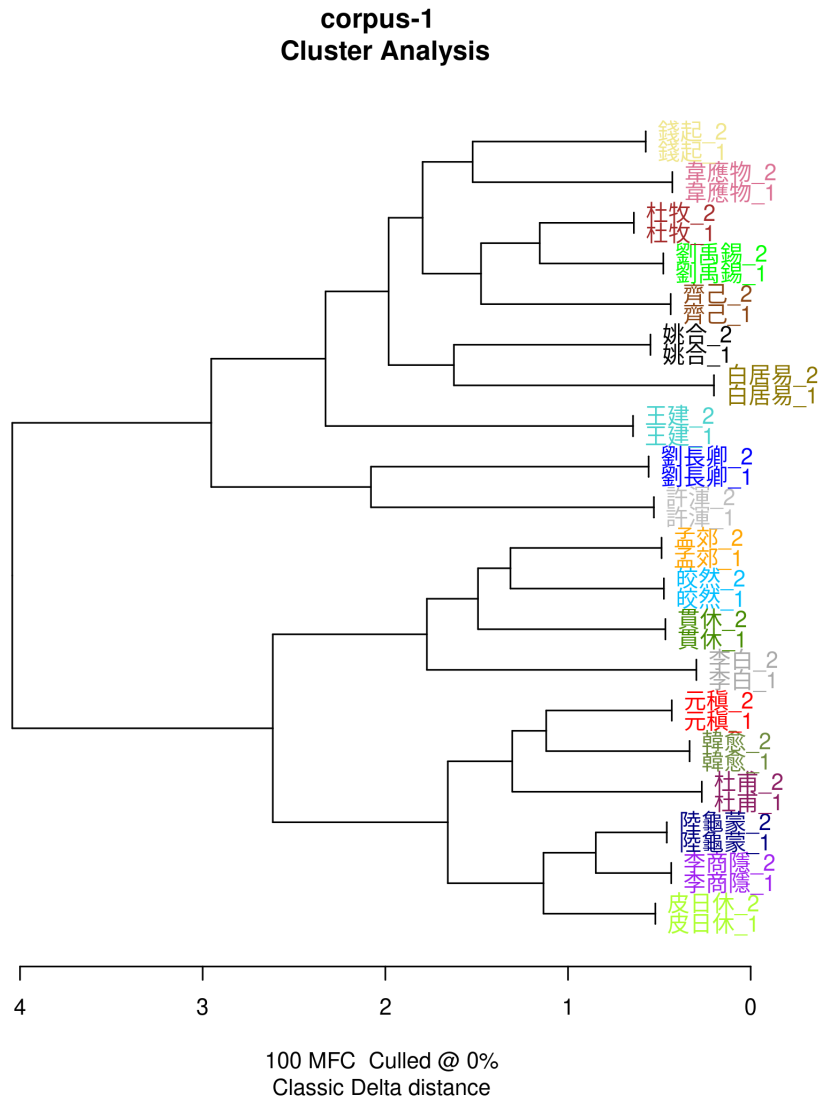


Figure 1: Cluster analysis of the first shuffled test corpus.

Table 1: Results: text distances

	元稹_1	元稹_2	劉禹錫_1	劉禹錫_2	劉長卿_1
元稹_1	0	0.4319	0.8917	0.8251	1.2953
元稹_2	0.4319	0	0.9412	0.8734	1.3814
劉禹錫_1	0.8917	0.9412	0	0.4782	1.1990
劉禹錫_2	0.8251	0.8734	0.4782	0	1.2040
劉長卿_1	1.2953	1.3814	1.1990	1.2040	0
劉長卿_2	1.4421	1.4649	1.3361	1.3059	0.5586



## Acknowledgements

I am grateful to Mariana Zorkina for her suggestions regarding the text corpus, specific bibliography points, and some observations about the research field, which we discussed in the early versions of this work. At the same time, all inaccuracies and errors in the text remain my responsibility.

## References

- John Burrows. 'Delta' : a measure of stylistic difference and a guide to likely authorship. *Literary and linguistic computing*, 17(3):267–287, 2002.
- Daniil Skorinkin and Boris Orekhov. Hacking stylometry with multiple voices: Imaginary writers can override authorial signal in delta. *Digital Scholarship in the Humanities*, 38(3):1247–1266, 2023.
- Maciej Eder, Jan Rybicki, and Mike Kestemont. Stylometry with R: a package for computational text analysis. *R Journal*, 8(1):107–121, 2016. URL <https://journal.r-project.org/archive/2016/RJ-2016-007/index.html>.
- David L Hoover. Testing burrows's delta. *Literary and linguistic computing*, 19(4):453–475, 2004.
- Antonio Miranda Garcia and Javier Calle Martin. Functionist words in authorship attribution studies. *Literary and Linguistic Computing*, 22(1):49–66, 2006.
- Fotis Jannidis and Gerhard Lauer. Burrows' s delta and its use in german literary history. *Distant readings. Topologies of German culture in the long nineteenth century*, pages 29–54, 2014.
- Laura Hernández-Lorenzo and Joanna Byszuk. Challenging stylometry: The authorship of the baroque play la segunda celestina. *Digital Scholarship in the Humanities*, 38(2):544–558, 2023.
- Jacques Savoy. Is starnone really the author behind ferrante? *Digital Scholarship in the Humanities*, 33(4):902–918, 2018.
- Jan Rybicki. Partners in life, partners in crime. *Drawing Elena Ferrante' s Profile*, pages 111–122, 2018.
- Jan Rybicki and Magda Heydel. The stylistics and stylometry of collaborative translation: Woolf' s night and day in polish. *Literary and Linguistic Computing*, 28(4):708–717, 2013.
- Daniil Skorinkin and Anastasiya Bonch-Osmolovskaya. "Osoby primety" v rechi khudozhestvennykh personazhej: kolichestvennyj analiz dialogov v "Voyne i mire" L. N. Tolstogo. *Istoriya: elektronnyy nauchno-obrazovatel'nyj zhurnal*, (7 (51)), 2016. doi:10.18254/S0001649-2-1.
- Ammar Adil AbdulRazzaq and Tareef Kamil Mustafa. Burrows-delta method fitness for arabic text authorship stylometric detection. *International Journal of Computer Science and Mobile Computing*, 3(6):69–78, 2014.
- Maciej Eder. Style-markers in authorship attribution: a cross-language study of the authorial fingerprint. *Studies in Polish Linguistics*, 6(1), 2011.
- Testing delta on chinese texts. In *Digital Humanities 2016: Conference Abstracts*, pages 781–783, Kraków, 2016. Jagiellonian University.
- Hang Xiao. On the applicability of zipf's law in chinese word frequency distribution. *J. Chin. Lang. Comput.*, 18: 33–46, 2008. URL <https://api.semanticscholar.org/CorpusID:8053003>.
- Chao-Lin Liu, Shuhua Zhang, Yuanli Geng, Hwei-ling Lai, and Hongsu Wang. Character distributions of classical chinese literary texts: Zipf's law, genres, and epochs. *arXiv preprint arXiv:1709.05587*, 2017.
- Shilei Huang and Jiangqin Wu. A pragmatic approach for classical chinese word segmentation. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- Donald Sturgeon. Unsupervised identification of text reuse in early chinese literature. *Digital Scholarship in the Humanities*, 33(3):670–684, 2018.
- snowtraces. poetry-source. <https://github.com/snowtraces/poetry-source>, 2020.
- Boris Orekhov. Delta for tang poets, Jul 2024. URL [osf.io/j4fn3](https://osf.io/j4fn3).