

# The State of Non-Traditional Authorship Attribution Studies – 2010: Some Problems and Solutions

**Rudman, Joseph**

jr20@heps.phys.cmu.edu  
Carnegie Mellon University, USA

In 1997, at the ACH-ALLC'97 conference at Queen's University, there was a session presented by R. Harald Baayen, David I. Holmes, Joe Rudman, and Fiona J. Tweedie, "The State of Authorship Attribution Studies: (1) The History and the Scope; (2) The Problems – Towards Credibility and Validity." Thirteen years have passed and well over 600 studies and papers dealing with non-traditional authorship attribution have been promulgated since that session.

This paper looks back at that session, a subsequent article published by Rudman in *Computers and the Humanities*, "The State of Authorship Attribution Studies: Some Problems and Solutions," and the more than 600 new publications. There are still major problems in the "science" on non-traditional authorship attribution. This paper goes on to assess the present state of the field – its successes, failures, and prospects.

## 1. Successes

It has been an exciting thirteen years with many advances. Each of the following (not a complete list) will be discussed:

1. Arguably, the most significant development in the field is the large contingent of computer scientists that have brought their perspectives to the table – led by Shlomo Argamon, Moshe Kopple, and a host of others.
2. The Dimacs Working Group on Developing Community.
3. Sir Brian Vickers' London Authorship Forum.
4. John Burrows' Busa Award.

5. Forensic Linguistics.
6. Successful studies such as Foster's *Primary Colors* work.
7. The continuing advances of practitioners such as John Burrows, David Hoover, Matthew Jockers, David Holmes, and others.
8. John Nerbonne's reissue of Mosteller and Wallace's *Applied Bayesian and Classical Inference: The Case of "The Federalist Papers."*
9. Patrick Juola's "Ad Hoc Authorship Attribution Competition." and his NSF funded JGAAP project.
10. The PAN Workshops. Uncovering Plagiarism, Authorship, and Social Software Misuse.

## 2. Acceptance

Contrary to what many practitioners of the non-traditional proclaim, there is not wide-spread acceptance of the field.

There have been many high profile problems with the concomitant negative publicity, e.g.:

1. Foster's misattribution of *A Funerall Elegie*
2. Foster's misattribution of the Jon Benét ransom note
3. Burrows' attribution then de-attribution of "A Vision"
4. The continuing bashing of Morton's CUSUM

Burrows' shift is something that every good scientist should do – search for errors or improvements in their experimental methodology and self correct.

## 3. Failures and Shortcomings

After thirteen years of increasing activity, there is still no consensus as to correct methodology or technique. Most authorship studies are still governed by expediency, e.g.:

- The texts are not the correct ones but they were available
- The controls are not complete but it would have taken too long to obtain the correct ones

The "umbrella" problem remains – most non-traditional authorship practitioners do not understand what constitutes a valid study.

Problems in the following areas will be explicated and solutions proposed:

- Knowledge of the Field (i.e. the Bibliography)
  - The fact that there have been so many authorship studies is good -- the fact that they have been published in over 90 different journals makes a complete literature search time consuming and difficult which is not good. To make things even more difficult, add to this the more than 14 books, 22 chapters in books, the 80 conference papers, the 10 reports, 22 dissertations, 9 newspaper articles, the 10 on-line self published papers, 4 encyclopedia entries.
- Reproducibility – verification
- The Experimental Plan
- The Primary Data – This is a major problem that is almost universally side-stepped.
- Style markers – Function words, n-grams, etc.
- Cross Validation – necessary but not sufficient
- The Control Groups – Genre, gender, time frame, etc.
- The Statistics – A range of techniques will be discussed – e.g. Neural Nets, SVM's, Sequence Kernals, Nave Bayes
- The Presentation – visualization

#### 4. Conclusion

In conclusion, there is a discussion of our role as gatekeepers:

- Rudman's caution that attribution studies on the *Historia Augusta* are an exercise in futility.
- Hoover and Argamon's modification and clarification of Burrows' Delta.
- Rudman's "Ripost" of Burrows' "History of Ophelia."
- Should we oppose patents such as Chaski's?
- The Daubert triangle.

---

#### References

**Argamon, Shlomo, et al.** (2003). 'Gender, Genre, and Writing Style in Formal Written Texts'. *Text*. **23.3**: 321-346.

**Baayen, Harald, Hans van Halteren, Anneke Neijt, and Fiona Tweedie.** (2002). 'An Experiment in Authorship Attribution'. *JADT 2002:6es Journées Internationales d'Analyse Statistique des Données Textuelles*.

**Brennan, Michael, and Rachel Greenstadt.** *Practical Attacks Against Authorship Attribution Techniques*. <http://www.cs.drexel.edu/greenie/brennan-paper.pdf> (accessed July 14, 2009).

**Burrows, John** (2007). 'Sarah and Henry Fielding and the Authorship of The History of Ophelia: A Computational Analysis'. *Script & Print*. **30.2**: 69-92.

**Chung, Cindy, and James PenneBaker** (2007). 'The Psychological Functions of Function Words'. *Social Communication*. K. Fiedler (ed.). New York: Psychology Press, pp. 343-359.

**Feiguina, Ol'ga, and Graeme Hirst** (2007). 'Authorship Attribution for Small Texts: Literary and Forensic Experiments'. *Proceedings of SIGIR '07 Workshop on Plagiarism Analysis, Authorship Identification, and Near-Duplicate Detection*. Amsterdam.

**Foster, Donald W.** (February 26, 1996). *Primary Culprit: An Analysis of a Novel of Politics*. New York, pp. 50-57.

**Khosmood, Foaad, and Robert Levinson** (2006). 'Toward Unification of Source Attribution Processes and Techniques'. *Proceedings of the Fifth International Conference on Machine Learning and Cybernetics*. Dalian, pp. 4551-4556.

**Love, Harold** (2002). *Attributing Authorship: An Introduction*. Cambridge: Cambridge University Press.

**Nieder Korn, William S.** (20 June 2002). *The New York Times*, B1, B5.

**Ramyaa, Congzhou, and Khaled Rasheed** (2004). 'Using Machine Learning Techniques for Stylometry'. *International Conference on Machine Learning (MLMTA'2004)*. Las Vegas.

**Rudman, Joseph** (2007). 'Sarah and Henry Fielding and the Authorship of "The History of Ophelia": A Ripost'. *Script & Print*. **31.3**: 147-163.

**Solon, Lawrence M., and Peter M. Tiersma** (2005). *Speaking of Crime*. Chicago: The University of Chicago Press.

**Stamatatos, Efstathios, Nikos Fakotakis, and George Kokkinakis** (2001). 'Automatic Text Categorization in Terms of Genre and Author'. *Computational Linguistics*. **26.4**: 471-495.

**Stein, Benno, et al. (eds.)** (2009). *PAN'09*.

**Tambouratzis, George** (2001). 'Assessing the Effectiveness of Feature Groups in Author Recognition Tasks with the SOM Model'. *IEEE Transactions on Systems, Man, and Cybernetics – Part C: Applications and Reviews*. **36.2**: 249-259.

**Van Halteren, Hans** (2004). 'Linguistic Profiling for Authorship Recognition and Verification'. *42nd Annual Meeting of the Association for Computational Linguistics*. Barcelona.