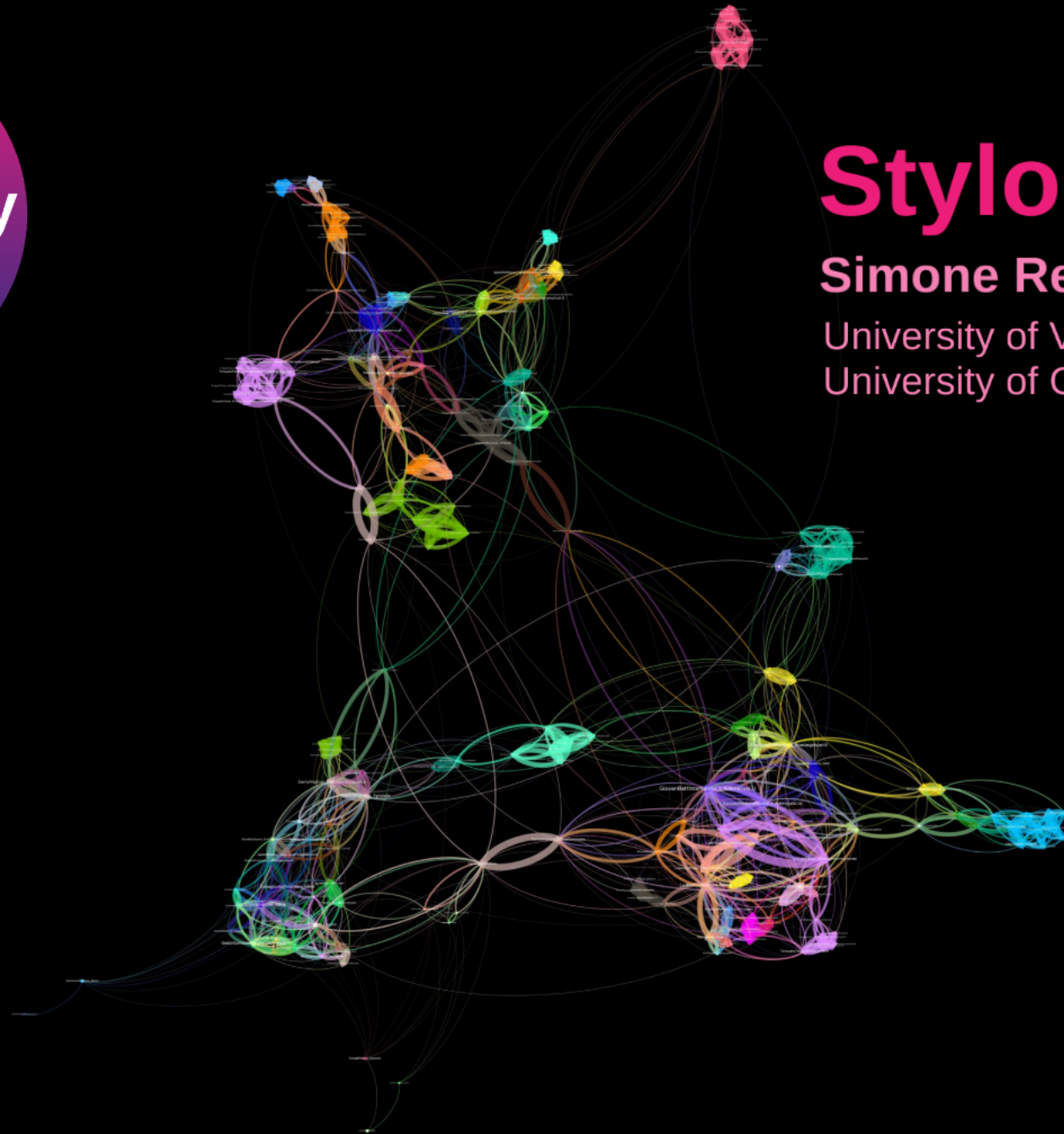**Stylometry**

**Authorship Attribution**

**Network Analysis**

# Stylometry

**Simone Rebora**

University of Verona
University of Göttingen

# Stylometry

# =

# "measuring style"

**The Origins**

**The Revolution**

# The Origins
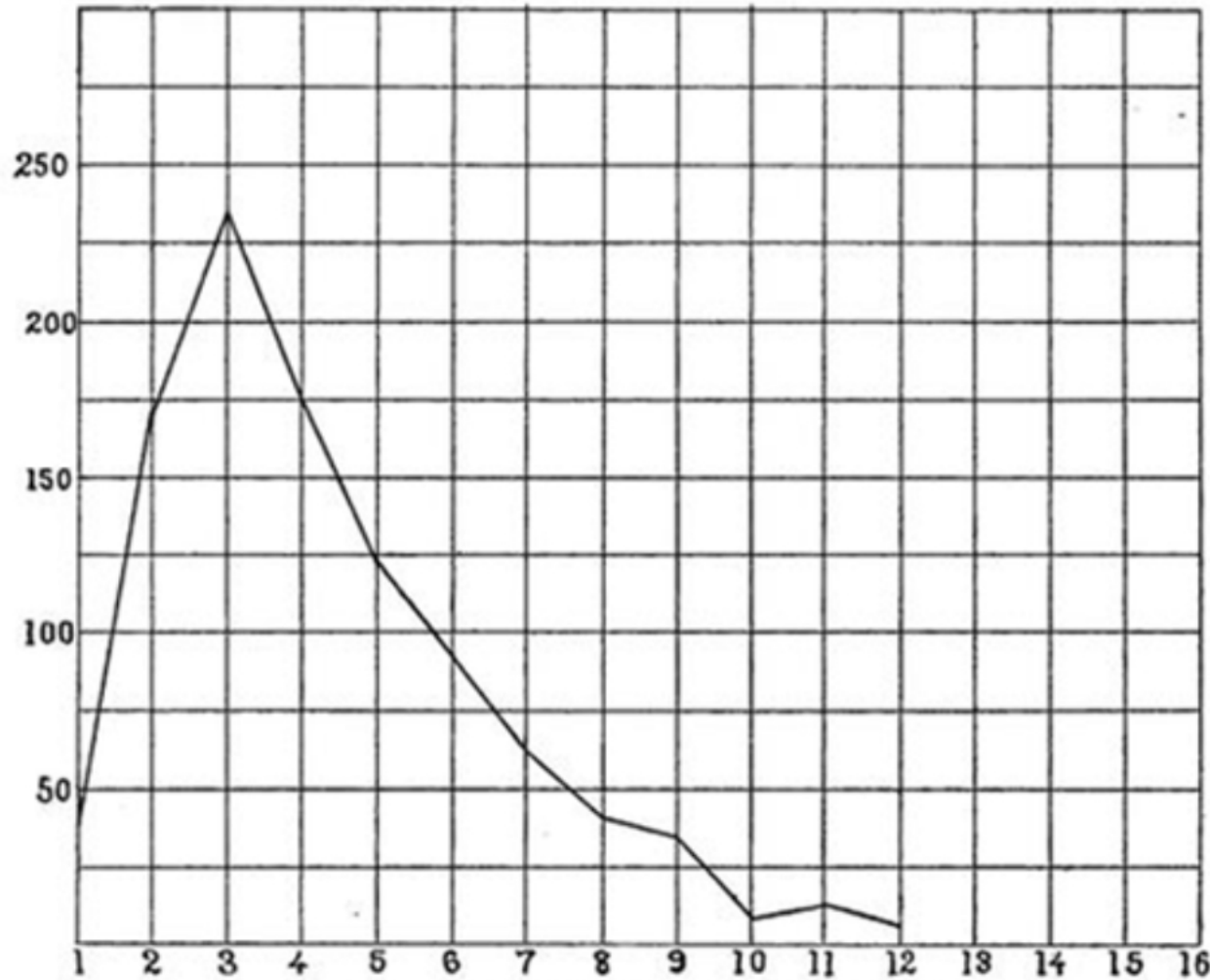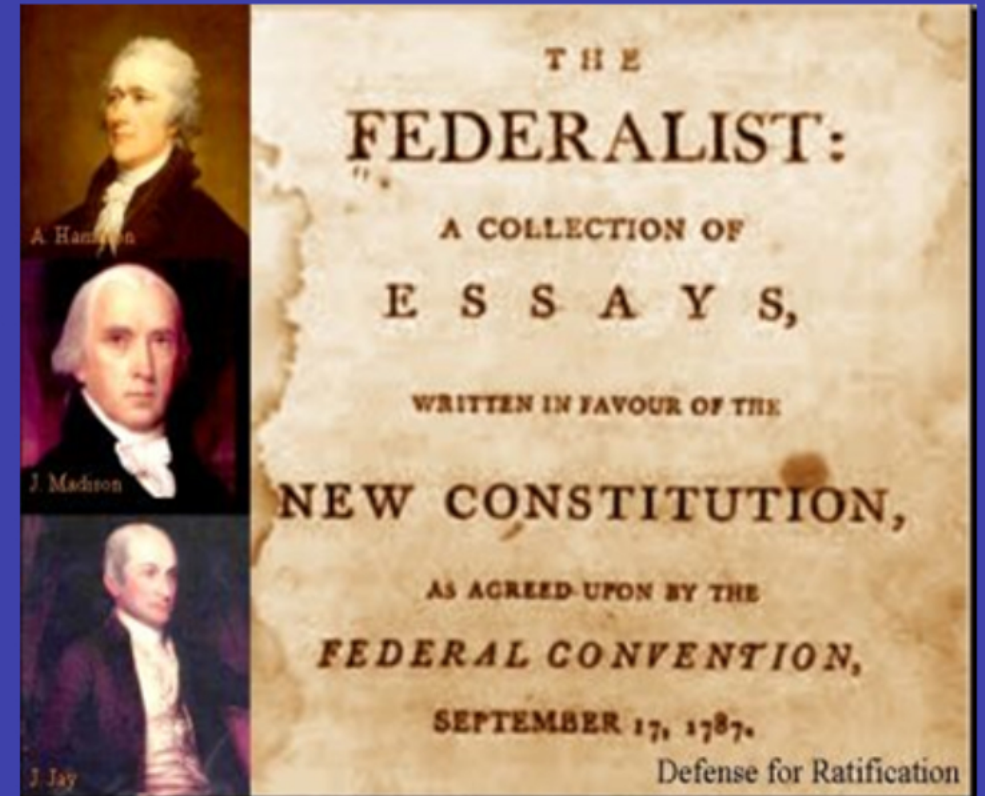


FIG. 1. — FIRST ONE THOUSAND WORDS IN 'OLIVER TWIST.'

Mendenhall, T. C. (1887). "The Characteristic Curves of Composition". Science. IX (214): 237–248

# A history of successes...

- 3 authors (A. Hamilton, J. Madison, J. Jay)
- 85 articles and essays written in 1787-1788, under the pseydonym "Publius"
- frequency of 165 words (mainly functional)
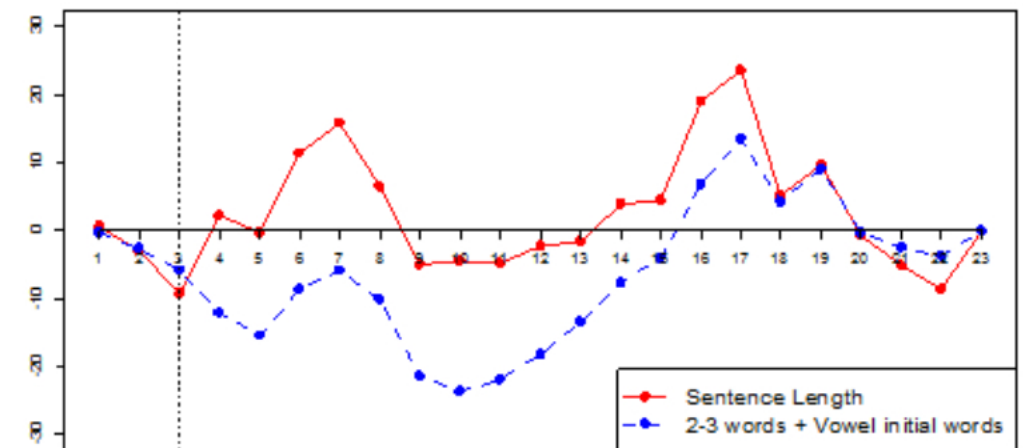
| | enough | while | whilst | upon |
|---|---|---|---|---|
| Hamilton | 0.59 | 0.26 | 0 | 2.93 |
| Madison | 0 | 0 | 0.47 | 0.16 |
| Disputed texts | 0 | 0 | 0.34 | 0.08 |
| Co-authored texts | 0.18 | 0 | 0.36 | 0.36 |

Mosteller & Wallace (1964)

# ...and Epic Failures

- Andrew Morton in the early '60 adapted Cumulative Sum – CUSUM or QSUM (a method which originally was used in the industrial quality control) to be used in texts.
- BBC live show (1993)

Documents of convicted criminals were attributed to … the Secretary of State for Justice!!!

# The Revolution

"Literary and Linguistic Computing" 17, no. 3 (2002): 267–87

## 'Delta': a Measure of Stylistic Difference and a Guide to Likely Authorship[1]

John Burrows
University of Newcastle, Australia

### Abstract

This paper is a companion to my 'Questions of authorship: attribution and beyond', in which I sketched a new way of using the relative frequencies of the very common words for comparing written texts and testing their likely authorship. The main emphasis of that paper was not on the new procedure but on the broader consequences of our increasing sophistication in making such comparisons and the increasing (although never absolute) reliability of our inferences about authorship. My present objects, accordingly, are to give a more complete account of the procedure itself; to report the outcome of an extensive set of trials; and to consider the strengths and limitations of the new procedure. The procedure offers a simple but comparatively accurate addition to our current methods of distinguishing the most likely author of texts exceeding about 1,500 words in length. It is of even greater value as a method of reducing the field of likely candidates for texts of as little as 100 words in length. Not unexpectedly, it
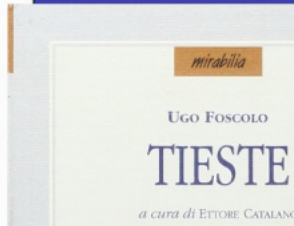
# Delta Distance

1. e
2. che
3. di
4. la
5. a
6. il
7. non
8. l
9. in
10. per
11. le
12. si
13. con
14. i
15. è
16. un
17. del
18. da
19. più
20. d
21. gli
22. ma

1. e
2. che
3. di
4. la
5. a
6. il
7. non
8. l
9. in
10. per
11. le
12. si
13. con
14. i
15. è
16. un
17. del
18. da
19. più
20. d
21. gli
22. ma

0.13
0.11
0.09
...

ALESSANDRO MAN...
ADELCHI

mirabilia
UGO FOSCOLO
TIESTE
a cura di ETTORE CATALANO

0.10
0.14
0.08
...

Alfieri
Tragedie

Goldoni
La locandiera
Introduzione di Giorgio Strehler

0.11
0.12
0.10
...

0.13
0.10
0.07
...

$$\Delta_{(AB)} = \frac{1}{n} \sum_{i=1}^{n} \left| \frac{A_i - \mu_i}{\sigma_i} - \frac{B_i - \mu_i}{\sigma_i} \right|$$

| | AlessandroManzoni_Adelchi | AlessandroManzoni_IlContediCarmagnola | AlessandroManzoni_InniSacri | AlessandroManzoni_Odi | AlessandroManzoni_Poesiegio |
|---|---|---|---|---|---|
| AlessandroManzoni_Adelchi | 0 | 0,481290655 | 0,666926925 | 0,738545533 | 0,568 |
| AlessandroManzoni_IlContediCarmagnola | 0,481290655 | 0 | 0,746348745 | 0,814261157 | 0,654 |
| AlessandroManzoni_InniSacri | 0,666926925 | 0,746348745 | 0 | 0,633663965 | 0,634 |
| AlessandroManzoni_Odi | 0,738545533 | 0,814261157 | 0,633663965 | 0 | 0,733 |
| AlessandroManzoni_Poesiegiovanili | 0,568820863 | 0,654375023 | 0,634854567 | 0,733827682 | |
| CarloGoldoni_GI'Innamorati | 0,980786338 | 0,936018177 | 1,013723738 | 1,101305203 | 0,950 |
| CarloGoldoni_IlCampiello | 1,016924762 | 1,031300757 | 1,018625104 | 1,092680684 | 0,929 |
| CarloGoldoni_IlServitorediduePadroni | 0,94860233 | 0,926662976 | 0,976288639 | 1,080804722 | 0,918 |
| CarloGoldoni_IlTeatrocomico | 0,915941412 | 0,896367382 | 0,971870697 | 1,085346366 | 0,898 |
| CarloGoldoni_IlVentaglio | 1,011953514 | 1,00041649 | 1,074888328 | 1,131792245 | 0,997 |
| CarloGoldoni_IRusteghi | 1,089096895 | 1,124315967 | 1,047451935 | 1,1240649 | 0,977 |
| CarloGoldoni_LaBottegadelcaffé | 0,997940632 | 0,980781404 | 1,069965126 | 1,139058754 | 0,993 |
| CarloGoldoni_LaFamigliadell'Antiquario | 0,97647637 | 0,968110166 | 1,038499373 | 1,080510085 | 0,953 |
| CarloGoldoni_LaLocandiera | 0,97946604 | 0,952399004 | 1,052505983 | 1,110322738 | 0,956 |
| CarloGoldoni_LeBaruffechiozzotte | 1,051753673 | 1,103993387 | 1,018834132 | 1,082447143 | 0,942 |
| CarloGoldoni_LeFemminepuntigliose | 0,940334542 | 0,938723973 | 1,008461186 | 1,076438004 | 0,917 |
| CarloGoldoni_LeSmanieperlaVilleggiatura | 1,023938091 | 0,964832878 | 1,056736183 | 1,148650567 | 1,007 |
| CarloGoldoni_UnadelleultimeserediCarnovale | 1,045847956 | 1,085480986 | 1,047945641 | 1,10681856 | 0,948 |
| VittorioAlfieri_Agamennone | 0,684514153 | 0,743793265 | 0,829452563 | 0,905939302 | 0,70 |
| VittorioAlfieri_Antigone | 0,73781244 | 0,801189414 | 0,824156384 | 0,91495815 | 0,721 |
| VittorioAlfieri_Brutosecondo | 0,675393312 | 0,675937144 | 0,830722082 | 0,910174086 | 0,668 |
| VittorioAlfieri_Filippo | 0,69672213 | 0,73856813 | 0,806194725 | 0,93419818 | 0,669 |
| VittorioAlfieri_MariaStuarda | 0,693145931 | 0,715015202 | 0,806081448 | 0,948928306 | 0,673 |
| VittorioAlfieri_Merope | 0,735463235 | 0,783055974 | 0,855979157 | 0,971583955 | 0,709 |
| VittorioAlfieri_Mirra | 0,76329317 | 0,819104452 | 0,864045202 | 0,9659327 | 0,760 |
| VittorioAlfieri_Oreste | 0,70530237 | 0,777981376 | 0,829335057 | 0,930970217 | 0,715 |
| VittorioAlfieri_Ottavia | 0,762895099 | 0,791949819 | 0,874379901 | 0,96265065 | 0,722 |
| VittorioAlfieri_Saul | 0,645417404 | 0,735038238 | 0,760393582 | 0,871007648 | 0,666 |

# Dendrograms

## Frequent Collocations and Authorial Style

David L. Hoover
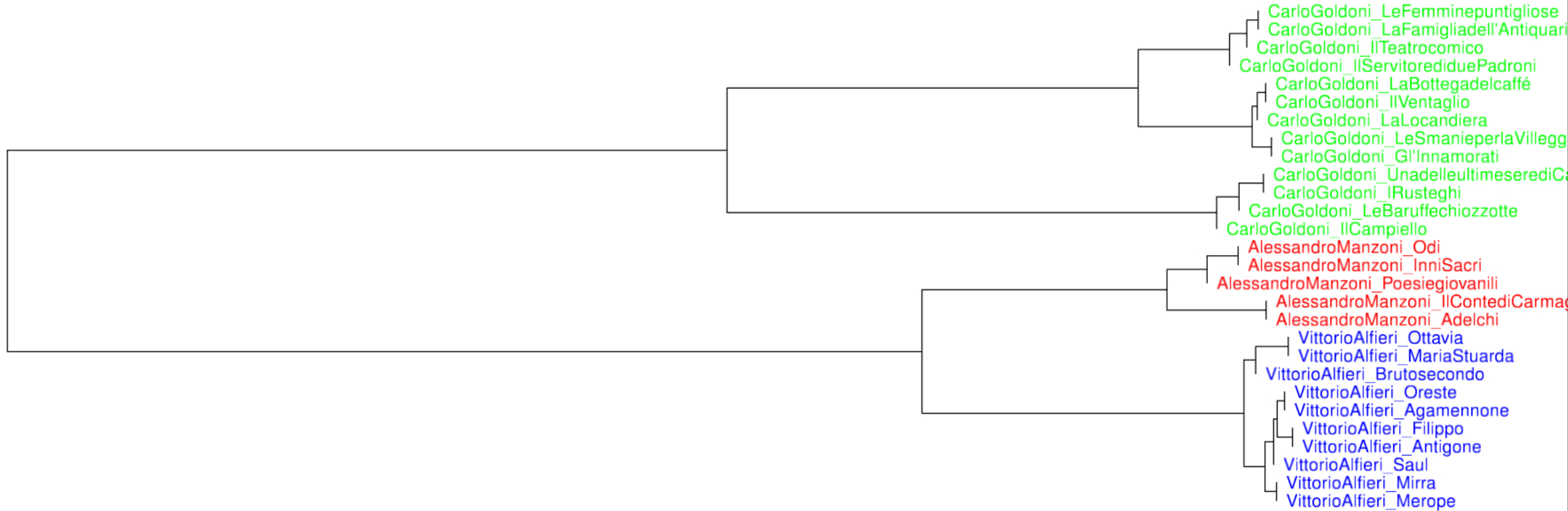New York University, New York, USA

### Abstract

This paper examines the effectiveness of multivariate analysis of the frequencies of frequent collocations in characterizing authorial style. Cluster analyses of collocations over various spans, types, and linkages are performed on groups of texts by known authors to determine how well the frequencies of those collocations correctly attribute the texts to their authors and distinguish them from texts by other authors. In each case the results are compared with those based on the frequencies of frequent words and the frequencies of frequent sequences of words. Cluster analyses based on frequent words and sequences ascribe many of the texts to their correct authors. However, analyses based on frequent collocations are more accurate for several groups of texts, sometimes producing more completely correct attributions than analyses based on either words or sequences and sometimes producing the only completely correct attributions. They also produce results for small groups of problematic novels and critical texts extracted from the larger corpora that are often superior to those based on

"Literary and Linguistic Computing" 18, no. 3 (2003): 261–83

# Letteratura Italiana
## Cluster Analysis

CarloGoldoni_LeFemminepuntigliose
CarloGoldoni_LaFamigliadell'Antiquari
CarloGoldoni_IlTeatrocomico
CarloGoldoni_IlServitorediduePadroni
CarloGoldoni_LaBottegadelcaffé
CarloGoldoni_IlVentaglio
CarloGoldoni_LaLocandiera
CarloGoldoni_LeSmanieperlaVillegg
CarloGoldoni_Gl'Innamorati
CarloGoldoni_UnadelleultimeserediCa
CarloGoldoni_IRusteghi
CarloGoldoni_LeBaruffechiozzotte
CarloGoldoni_IlCampiello
AlessandroManzoni_Odi
AlessandroManzoni_InniSacri
AlessandroManzoni_Poesiegiovanili
AlessandroManzoni_IlContediCarmag
AlessandroManzoni_Adelchi
VittorioAlfieri_Ottavia
VittorioAlfieri_MariaStuarda
VittorioAlfieri_Brutosecondo
VittorioAlfieri_Oreste
VittorioAlfieri_Agamennone
VittorioAlfieri_Filippo
VittorioAlfieri_Antigone
VittorioAlfieri_Saul
VittorioAlfieri_Mirra
VittorioAlfieri_Merope

8        6        4        2        0

100 MFW  Culled @ 0%
Classic Delta distance

Authorship Attribution

Juola vs. Rowling

Musil Project (in Verona)

# Juola vs. Rowling



Photo by Daniel Ogren

# Musil Project (in Verona)

# Letteratura Italiana
# Cluster Analysis

CarloGoldoni_LeFemminepuntigliose
CarloGoldoni_LaFamigliadell'Antiquari
CarloGoldoni_IlTeatrocomico
CarloGoldoni_IlServitorediduePadroni
CarloGoldoni_LaBottegadelcaffé
CarloGoldoni_IlVentaglio
CarloGoldoni_LaLocandiera
CarloGoldoni_LeSmanieperlaVillegg
CarloGoldoni_Gl'Innamorati
CarloGoldoni_UnadelleultimeserediC
CarloGoldoni_IRusteghi
CarloGoldoni_LeBaruffechiozzotte
CarloGoldoni_IlCampiello
AlessandroManzoni_Odi
AlessandroManzoni_InniSacri
AlessandroManzoni_Poesiegiovanili
AlessandroManzoni_IlContediCarmag
AlessandroManzoni_Adelchi
VittorioAlfieri_Ottavia
VittorioAlfieri_MariaStuarda
VittorioAlfieri_Brutosecondo
VittorioAlfieri_Oreste
VittorioAlfieri_Agamennone
VittorioAlfieri_Filippo
VittorioAlfieri_Antigone
VittorioAlfieri_Saul
VittorioAlfieri_Mirra
VittorioAlfieri_Merope

8    6    4    2    0

100 MFW  Culled @ 0%
Classic Delta distance

# Letteratura Italiana
## Cluster Analysis



LudovicoAriosto_Satire
LudovicoAriosto_Rime
LudovicoAriosto_OrlandoFurio
LudovicoAriosto_OrlandoFurio
LudovicoAriosto_Icinquecanti
TorquatoTasso_IlReTorrismondo
TorquatoTasso_Aminta
TorquatoTasso_Rinaldo
TorquatoTasso_LaGerusalemmeLib
AlessandroManzoni_IlContediCarmagi
AlessandroManzoni_Adelchi
UgoFoscolo_Aiace
AlessandroManzoni_Odi
AlessandroManzoni_InniSacri
MatteoMariaBoiardo_OrlandoInn
MatteoMariaBoiardo_OrlandoInn
MatteoMariaBoiardo_OrlandoInnam
VittorioAlfieri_Brutosecondo
UgoFoscolo_Tieste
VittorioAlfieri_Mirra
VittorioAlfieri_Antigone
VittorioAlfieri_Ottavia
VittorioAlfieri_MariaStuarda
VittorioAlfieri_Merope
VittorioAlfieri_Filippo
VittorioAlfieri_Oreste
VittorioAlfieri_Agamennone
VittorioAlfieri_Saul
CarloGoldoni_LeFemminepuntiglio
CarloGoldoni_LaFamigliadell Antiq
CarloGoldoni_IlTeatrocomico
CarloGoldoni_IlServitoredidue Padroni
CarloGoldoni_LaLocandiera
CarloGoldoni_IlVentaglio
CarloGoldoni_LaBottegadelcaffè
CarloGoldoni_LeSmanieperlaVilleg
CarloGoldoni_Gl Innamorati
CarloGoldoni_Unadelleultimeseredi
CarloGoldoni_LeBaruffechiozzotte
CarloGoldoni_IRusteghi
CarloGoldoni_IlCampiello

8    6    4    2    0

100 MFW  Culled @ 0%
Classic Delta distance

LudovicoAriosto_OrlandoFurio
LudovicoAriosto_OrlandoFurio
LudovicoAriosto_Icinquecanti
TorquatoTasso_IlReTorrismondo
TorquatoTasso_Aminta
TorquatoTasso_Rinaldo
TorquatoTasso_LaGerusalemmeLib
AlessandroManzoni_IlContediCarmag
AlessandroManzoni_Adelchi
UgoFoscolo_Aiace
AlessandroManzoni_Odi
AlessandroManzoni_InniSacri
MatteoMariaBoiardo_OrlandoInn
MatteoMariaBoiardo_OrlandoInn
MatteoMariaBoiardo_OrlandoInnam
VittorioAlfieri_Brutosecondo
UgoFoscolo_Tieste
VittorioAlfieri_Mirra
VittorioAlfieri_Antigone
VittorioAlfieri_Ottavia

**Letteratura Italiana**
**Cluster Analysis**

LudovicoAriosto_Satire
LudovicoAriosto_Rime
LudovicoAriosto_OrlandoFurios
LudovicoAriosto_OrlandoFurios
LudovicoAriosto_Icinquecanti
TorquatoTasso_IlReTorrismondo
TorquatoTasso_Aminta
TorquatoTasso_Rinaldo
TorquatoTasso_LaGerusalemmeLiber
MatteoMariaBoiardo_OrlandoInnam
MatteoMariaBoiardo_OrlandoInnam
MatteoMariaBoiardo_OrlandoInnamo
UgoFoscolo_Tieste
UgoFoscolo_Ajace
AlessandroManzoni_IlContediCarmagn
AlessandroManzoni_Adelchi
AlessandroManzoni_Odi_
AlessandroManzoni_InniSacri
VittorioAlfieri_Ottavia
VittorioAlfieri_MariaStuarda
VittorioAlfieri_Mirra
VittorioAlfieri_Brutosecondo
VittorioAlfieri_Oreste
VittorioAlfieri_Agamennone
VittorioAlfieri_Filippo
VittorioAlfieri_Merobe
VittorioAlfieri_Antigone
VittorioAlfieri_Saul
CarloGoldoni_LeFemminepuntigliose
CarloGoldoni_LaFamigliadell'Antiqua
CarloGoldoni_IlTeatrocomico
CarloGoldoni_IlServitorediduePadroni
CarloGoldoni_IlVentaglio
CarloGoldoni_GI'Innamorati
CarloGoldoni_LaLocandiera
CarloGoldoni_LaBottegadelcaffè
CarloGoldoni_LeSmanieperlaVilleggia
CarloGoldoni_Unadelleultimeserediq
CarloGoldoni_IRusteghi
CarloGoldoni_LeBaruffechiozzotte
CarloGoldoni_IlCampiello

8          6          4          2          0

200 MFW  Culled @ 0%
Classic Delta distance

# Letteratura Italiana
## Cluster Analysis



TorquatoTasso_IlReTorrismondo
TorquatoTasso_Aminta
TorquatoTasso_Rinaldo
TorquatoTasso_LaGerusalemmeLibera
LudovicoAriosto_Satire
LudovicoAriosto_Rime
LudovicoAriosto_OrlandoFurioso
LudovicoAriosto_OrlandoFurioso
LudovicoAriosto_Icinquecanti
AlessandroManzoni_Odi
AlessandroManzoni_InniSacri
AlessandroManzoni_IlContediCarmagno
AlessandroManzoni_Adelchi
MatteoMariaBoiardo_OrlandoInnam
MatteoMariaBoiardo_OrlandoInnam
MatteoMariaBoiardo_OrlandoInnamo
VittorioAlfieri_Oreste
VittorioAlfieri_Agamennone
VittorioAlfieri_Filippo
VittorioAlfieri_Antigone
VittorioAlfieri_Merope
VittorioAlfieri_Ottavia
VittorioAlfieri_MariaStuarda
VittorioAlfieri_Mirra
VittorioAlfieri_Brutosecondo
VittorioAlfieri_Saul
UgoFoscolo_Tieste
UgoFoscolo_Aiace
CarloGoldoni_LeFemminepuntigliose
CarloGoldoni_IlTeatrocomico
CarloGoldoni_LaFamigliadell_Antiquario
CarloGoldoni_IlServitorediduePadroni
CarloGoldoni_IlVentaglio
CarloGoldoni_Gl_Innamorati
CarloGoldoni_LaLocandiera
CarloGoldoni_LaBottegadelcaffé
CarloGoldoni_LeSmanieperlaVilleggiatur
CarloGoldoni_Unadelleultimeseredic
CarloGoldoni_IRusteghi
CarloGoldoni_LeBaruffechiozzotte
CarloGoldoni_IlCampiello

6    4    2    0

300 MFW  Culled @ 0%
Classic Delta distance

# Stylometry with R: A Package for Computational Text Analysis

*by Maciej Eder, Jan Rybicki and Mike Kestemont*

**Abstract**  This software paper describes 'Stylometry with R' (stylo), a flexible R package for the high-level analysis of writing style in stylometry. Stylometry (computational stylistics) is concerned with the quantitative study of writing style, e.g. authorship verification, an application which has considerable potential in forensic contexts, as well as historical research. In this paper we introduce the possibilities of **stylo** for computational text analysis, via a number of dummy case studies from English and French literature. We demonstrate how the package is particularly useful in the exploratory statistical analysis of texts, e.g. with respect to authorial writing style. Because **stylo** provides an attractive graphical user interface for high-level exploratory analyses, it is especially suited for an audience of novices, without programming skills (e.g. from the Digital Humanities). More experienced users can benefit from our implementation of a series of standard pipelines for text processing, as well as a number of similarity metrics.

## Introduction

Authorship is a topic which continues to attract considerable attention with the larger public. This claim is well illustrated by a number of high-profile case studies that have recently made headlines across the popular media, such as the attribution of a pseudonymously published work to acclaimed
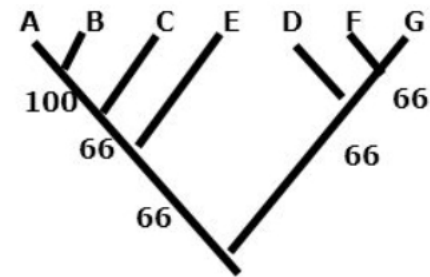
# Consensus Trees

## Majority rule consensus


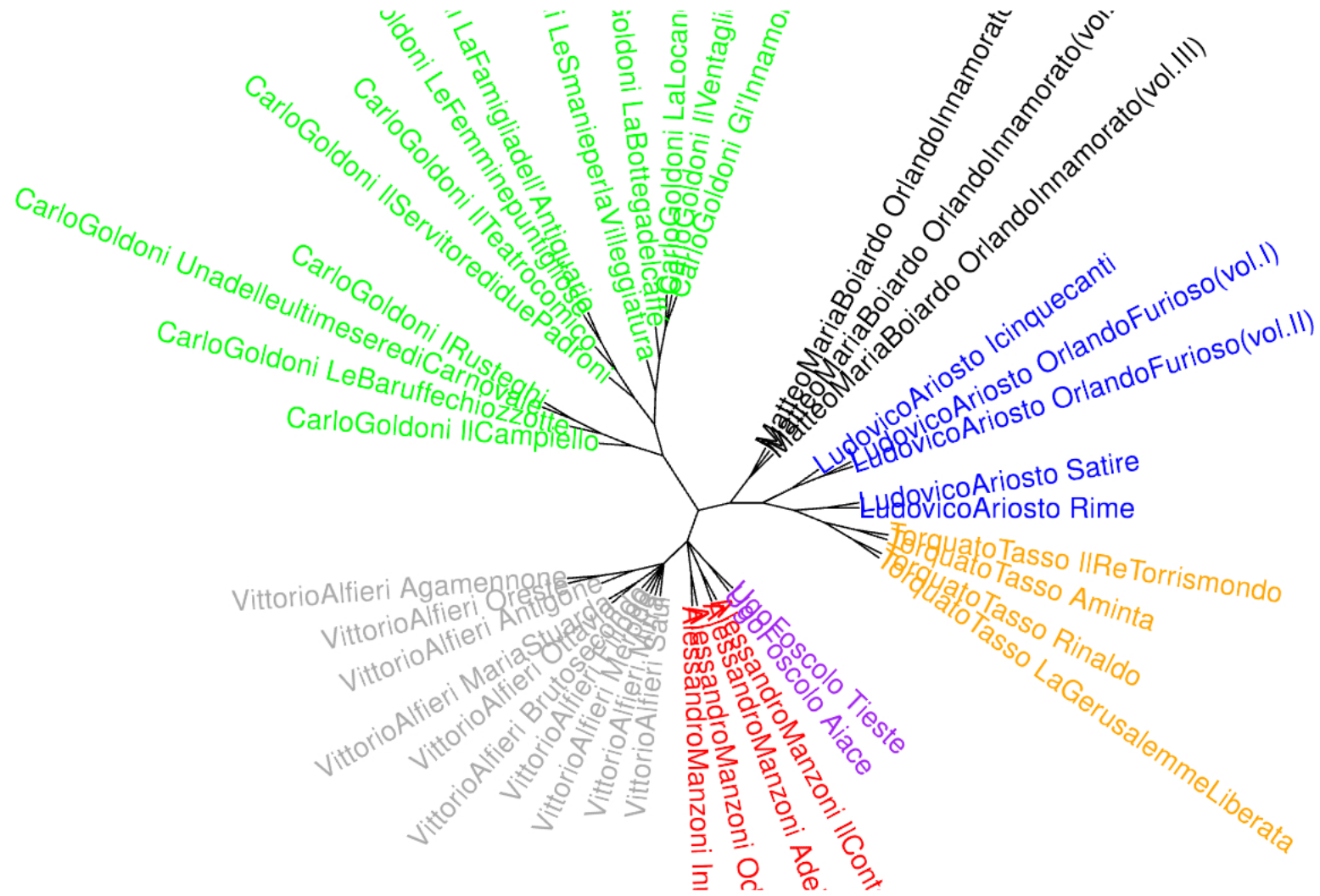
Numbers indicate frequency of clades in the fundamental trees
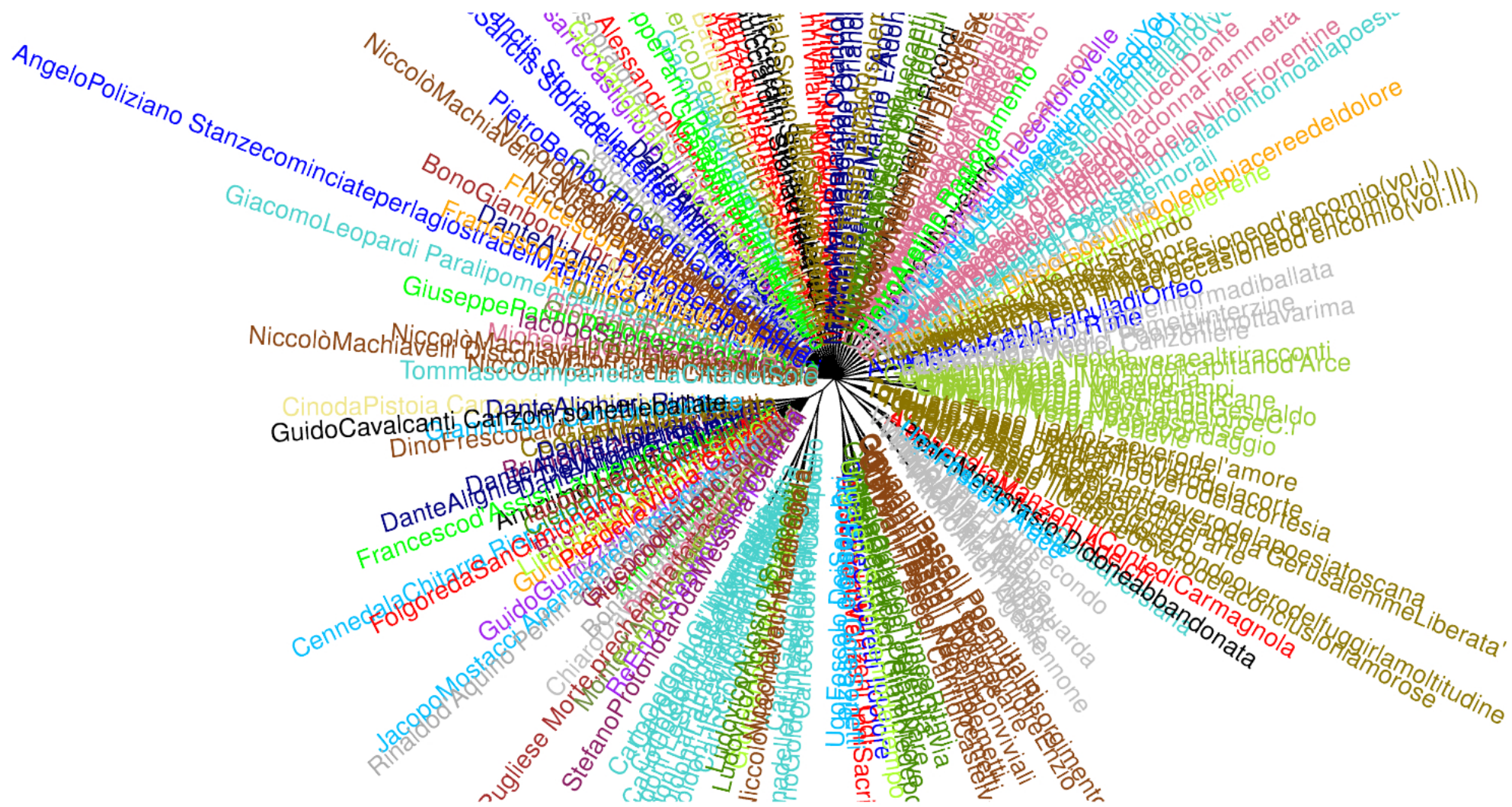
**MAJORITY-RULE CONSENSUS TREE**

# Letteratura Italiana
## Bootstrap Consensus Tree



CarloGoldoni UnadelleultimeserediCarnovale
CarloGoldoni LeBaruffechiozzotte
CarloGoldoni IlCampiello
CarloGoldoni IlServitoredeiduePadroni
CarloGoldoni IlTeatrocomico
CarloGoldoni IRusteghi
CarloGoldoni LeFemminepuntigliose
CarloGoldoni LaFamigliadell'Antiquario
CarloGoldoni LeSmanieperlaVilleggiatura
CarloGoldoni LaBottegadelcaffe
CarloGoldoni LaLocan...
CarloGoldoni IlVentagli...
CarloGoldoni Gl'Innamo...

MatteoMariaBoiardo OrlandoInnamorat...
MatteoMariaBoiardo OrlandoInnamorato(vol...
MatteoMariaBoiardo OrlandoInnamorato(vol.III)

LudovicoAriosto Icinquecanti
LudovicoAriosto OrlandoFurioso(vol.I)
LudovicoAriosto OrlandoFurioso(vol.II)
LudovicoAriosto OrlandoFurioso(vol.III)
LudovicoAriosto Satire
LudovicoAriosto Rime

TorquatoTasso IlReTorrismondo
TorquatoTasso Aminta
TorquatoTasso Rinaldo
TorquatoTasso LaGerusalemmeLiberata

VittorioAlfieri Agamennone
VittorioAlfieri Oreste
VittorioAlfieri Antigone
VittorioAlfieri MariaStuarda
VittorioAlfieri Ottavia
VittorioAlfieri Brutosecondo
VittorioAlfieri Merope
VittorioAlfieri Saul

UgoFoscolo Tieste
UgoFoscolo Aiace

AlessandroManzoni IlCont...
AlessandroManzoni Ade...
AlessandroManzoni Oc...
AlessandroManzoni In...

100−1000 MFW  Culled @ 0%
Classic Delta distance Consensus 0.5

# Letteratura Italiana
## Bootstrap Consensus Tree



100–1000 MFW  Culled @ 0%
Classic Delta distance Consensus 0.5

# Visualization in stylometry: Cluster analysis using networks

**Maciej Eder**

Pedagogical University of Kraków, Poland

Institute of Polish Language, PAS

**Correspondence:**

Maciej Eder, Institute of Polish Studies, Pedagogical University of Kraków, ul. Podchorążych 2, 30-084 Kraków, Poland.

E-mail: maciejjeder@gmail.com

## Abstract

The aim of this article is to discuss reliability issues of a few visual techniques used in stylometry, and to introduce a new method that enhances the explanatory power of visualization with a procedure of validation inspired by advanced statistical methods. A promising way of extending cluster analysis dendrograms with a self-validating procedure involves producing numerous particular 'snapshots', or dendrograms produced using different input parameters, and combining them all into the form of a consensus tree. Significantly better results, however, can be obtained using a new visualization technique, which combines the idea of nearest neighborhood derived from cluster analysis, the idea of hammering out a clustering consensus from bootstrap consensus trees, with the idea of mapping textual similarities onto a form of a network. Additionally, network analysis seems to be a good solution for large data sets.

## 1 Introduction

Most of the computational methods used in stylometry have been originally introduced to solve authorship attribution problems. This fact had an algorithms, suitable for classification tasks, derived mostly from the field of biometrics, nuclear physics, or software engineering, that could be easily adopted to authorship attribution. They include naïve Bayes classification, support vector machines

Fig. 6. Two algorithms of mapping textual relations: establishing weighted links to a nearest neighbor and two runners-up (top); producing a consensus network (bottom).

Casa_Galateo

PietroAretino_Dialogo

PietroAretino_Ragionamento

BonoGianboni_Libro de'viziedellevirtudi

GiovanniBoccaccio_TrattatelloinlaudediDante

FrancoSacchetti_IlTrecentonovelle

GiovanniBoccaccio_Corbaccio

MarcoPolo_IlMilione

GiovanniBoccaccio_Decameron

Anonimo_IlNovellino

GiovanniBoccaccio_ElegiadiMadonnaFiammetta

GiovanniBoccaccio_IlFilocolo