

Authorship Analysis based on Data Compression

Daniele Cerra, Mihai Datcu, and Peter Reinartz

German Aerospace Center (DLR), Muenchner str. 20, 82234 Wessling, Germany

Corresponding author's email: daniele.cerra@dlr.de,

phone: +49 8153 28-1496, fax: +49 8153 28-1444.

Abstract

This paper proposes to perform authorship analysis using the Fast Compression Distance (FCD), a similarity measure based on compression with dictionaries directly extracted from the written texts. The FCD computes a similarity between two documents through an effective binary search on the intersection set between the two related dictionaries. In the reported experiments the proposed method is applied to documents which are heterogeneous in style, written in five different languages and coming from different historical periods. Results are comparable to the state of the art and outperform traditional compression-based methods.

Keywords: Authorship Analysis, Data Compression, Similarity Measure

1. Introduction

The task of automatically recognizing the author of a given text finds several uses in practical applications, ranging from authorship attribution to plagiarism detection, and it is a challenging one (Stamatatos, 2009). While the structure of a document can be easily interpreted by a machine, the

13 description of the style of each author is in general subjective, and therefore
14 hard to derive in natural language; it is even harder to find a description
15 which enables a machine to automatically tell one author from the other. A
16 literature review on modern authorship attribution methods, usually coming
17 from the fields of machine learning and statistical analysis, is reported in
18 Stamatatos (2009); Jockers and Witten (2010); Koppel et al. (2009); Grieve
19 (2007); Juola (2006). Among these, algorithms based on similarity measures
20 such as Benedetto et al. (2002) and Koppel et al. (2011) are widely employed
21 and usually assign an anonymous text to the author of the most similar
22 document in the training data.

23 During the last decade, compression-based distance measures have been
24 effectively applied to cluster texts written by different authors (Cilibrasi and Vitányi,
25 2005) and to perform plagiarism detection (Chen et al., 2004). Such univer-
26 sal similarity measures, of which the most well-known is the Normalized
27 Compression Distance (NCD), employ general compressors to estimate the
28 amount of shared information between two objects. Similar concepts are
29 also used by methods using runlength histograms to retrieve and classify
30 documents (Gordo et al., 2013). Experiments carried out in Oliveira et al.
31 (2013) conclude that NCD-based methods for authorship analysis outper-
32 form state-of-the-art classification methodologies such as Support Vector
33 Machines. A study on larger and more statistically meaningful datasets
34 shows NCD-methods to be competitive with respect to the state of the art
35 (de Graaff, 2012), while Stamatatos (2009) reports that compression-based

36 methods are effective but hard to use in practice as they are very slow.

37 Indeed the universality of these measures comes at a price, as the com-
38 pression algorithm must be run at least n^2 times on n objects to derive a
39 distance matrix, slowing down the analysis. Furthermore, as these methods
40 are applied to raw data they cannot be tuned to increase their performance
41 on a given data type. We propose then to perform these tasks using the Fast
42 Compression Distance (FCD) recently defined in Cerra and Datcu (2012),
43 which provides superior performances with a reduced computational com-
44 plexity with respect to the NCD, and can be tuned according to the kind
45 of data at hand. In the case of natural texts, only FCD's general settings
46 should be adjusted according to the language of the dataset, thus keeping
47 the desirable parameter-free approach typical of NCD. Applications to au-
48 thorship and plagiarism analysis are derived by extracting meaningful dictio-
49 naries directly from the strings representing the data instances and matching
50 them. The reported experiments show that improvements over traditional
51 compression-based analysis can be dramatic, and that the FCD could be-
52 come an important tool of easy usage for the automated analysis of texts, as
53 satisfactory results are achieved skipping any parameters setting step. The
54 only exception is an optional text preprocessing step which only needs to
55 be set once for documents of a given language, and does not depend on the
56 specific dataset.

57 The paper is structured as follows. Section 2 introduces compression-
58 based similarity measures and the FCD, which will be validated in an array

59 of experiments reported in Section 3. We conclude in Section 4.

60 **2. Fast Compression Distance**

61 Compression-based similarity measures exploit general off-the-shelf com-
62 pressors to estimate the amount of information shared by any two objects.
63 They have been employed for clustering and classification on diverse data
64 types such as texts and images (Watanabe et al., 2002), with Keogh et al.
65 (2004) reporting that they outperform general distance measures. The most
66 widely known and used of such notions is the Normalized Compression Dis-
67 tance (NCD), defined for any two objects x and y as:

$$NCD(x, y) = \frac{C(x, y) - \min C(x), C(y)}{\max C(x), C(y)} \quad (1)$$

68 where $C(x)$ represents the size of x after being compressed by a com-
69 pressor (such as Gzip), and $C(x, y)$ is the size of the compressed version
70 of x appended to y . If $x = y$, the NCD is approximately 0, as the full
71 string y can be described in terms of previous strings found in x ; if x and y
72 share no common information the NCD is $1 + e$, where e is a small quantity
73 (usually $e < 0.1$) due to imperfections characterizing real compressors. The
74 idea is that if x and y share common information they will compress better
75 together than separately, as the compressor will be able to reuse recurring
76 patterns found in one of them to more efficiently compress the other. The
77 generality of NCD allows applying it to diverse datatypes, including natu-
78 ral texts. Applications to authorship categorization have been presented by

79 Cilibrasi and Vitányi (2005), while plagiarism detection of students assign-
80 ments has been successfully carried out by Chen et al. (2004).

81 A modified version of NCD based on the extraction of dictionaries has
82 been first defined by Macedonas et al. (2008). The advantages of using
83 dictionary-based methods have been then studied by Cerra and Datcu (2012),
84 in which the authors define a Fast Compression Distance (FCD), and succes-
85 fully apply it to image analysis. The algorithm can be used for texts analysis
86 as follows.

87 First of all, all special characters such as punctuation marks are removed
88 from a string x , which is subsequently tokenized in a set of words W_x . The
89 sequence of tokens is analysed by the encoding algorithm of the Lempel-
90 Ziv-Welch (LZW) compressor (Welch, 1984), with the difference that words
91 rather than characters are taken into account. The algorithm initializes the
92 dictionary $D(x)$ with all the words W_x . Then the string x is scanned for
93 successively longer sequences of words in $D(x)$ until a mismatch in $D(x)$ takes
94 place; at this point the code for the longest pattern p in the dictionary is sent
95 to output, and the new string (p + the last word which caused a mismatch)
96 is added to $D(x)$. The last input word is then used as the next starting
97 point: in this way, successively longer sequences of words are registered in
98 the dictionary and made available for subsequent encoding, with no repeated
99 entries in $D(x)$. An example for the encoding of the string "TO BE OR
100 NOT TO BE OR NOT TO BE OR WHAT" after tokenization is reported
101 in Table 1. It helps to remark that the output of the simulated compression

Table 1: LZW encoding of the tokens composing the string "TO BE OR NOT TO BE OR NOT TO BE OR WHAT". The compressor tries to substitute pattern codes referring to sequences of words which occurred previously in the text.

Current token	Next token	Output	Added to Dictionary
Null	TO		
TO	BE	<i>TO</i>	TO BE= $\langle 1 \rangle$
BE	OR	<i>BE</i>	BE OR= $\langle 2 \rangle$
OR	NOT	<i>OR</i>	OR NOT= $\langle 3 \rangle$
NOT	TO	<i>NOT</i>	NOT TO= $\langle 4 \rangle$
TO BE	OR	$\langle 1 \rangle$	TO BE OR= $\langle 5 \rangle$
OR NOT	TO	$\langle 3 \rangle$	OR NOT TO= $\langle 6 \rangle$
TO BE OR	WHAT	$\langle 5 \rangle$	TO BE OR WHAT= $\langle 7 \rangle$
WHAT	#	<i>WHAT</i>	

102 process is not of interest for us, as the only thing that will be used is the
 103 dictionary.

104 The patterns contained in the dictionary $D(x)$ are then sorted in ascend-
 105 ing alphabetical order to enable the binary search of each pattern in time
 106 $O(\log N)$, where N is the number of entries in $D(x)$. The dictionary is finally
 107 stored for future use: this procedure may be carried out offline and has to be
 108 performed only once for each data instance. Whenever a string x is checked
 109 against a database containing n dictionaries, a dictionary $D(x)$ is extracted
 110 from x as described and matched against each of the n dictionaries. The
 111 FCD between x and an object y represented by $D(y)$ is defined as:

$$FCD(x, y) = \frac{|D(x)| - \cap(D(x), D(y))}{|D(x)|} \quad (2)$$

112 where $|D(x)|$ and $|D(y)|$ are the sizes of the relative dictionaries, repre-
 113 sented by the number of entries they contain, and $\cap(D(x), D(y))$ is the num-

114 ber of patterns which are found in both dictionaries. We have $FCD(x, y) = 0$
115 iff all patterns in $D(x)$ are contained also in $D(y)$, and $FCD(x, y) = 1$ if no
116 single pattern is shared between the two objects.

117 The FCD allows computing a compression-based distance between two
118 objects in a faster way with respect to NCD (up to one order of magnitude),
119 as the dictionary for each object must be extracted only once and comput-
120 ing the intersection between two dictionaries $D(x)$ and $D(y)$ is faster than
121 compressing the concatenation of x appended to y (Cerra and Datcu, 2012).
122 The FCD is also more accurate, as it overcomes drawbacks such as the lim-
123 ited size of the lookup tables, which are employed by real compressors for
124 efficiency constraints: this allows exploiting all the patterns contained in a
125 string. Furthermore, while the NCD is totally data-driven, the FCD enables
126 a token-based analysis which allows preprocessing the data, by decompos-
127 ing the objects into fragments which are semantically relevant for a given
128 data type or application. This constitutes a great advantage in the case of
129 plain texts, as the direct analysis of words contained in a document and
130 their concatenations allows focusing on the relevant informational content.
131 In plain English, this means that the matching of substrings in words which
132 may have no semantic relation between them (e.g. ‘butter’ and ‘butterfly’)
133 is prevented. Additional improvements can be made depending on the texts
134 language. For the case of English texts, the suffix ‘s’ can be removed from
135 each token, while from documents in Italian it helps to remove the last vowel
136 from each word: this avoids considering semantically different plurals and

137 some verbal forms.

138 A drawback of the proposed method is that it cannot be applied effectively
139 to very short texts. The algorithm needs to find reoccurring word sequences
140 in order to extract dictionaries of a relevant size, which are needed in order
141 to find patterns shared with other dictionaries. Therefore, the compression
142 of the initial part of a string is not effective: we estimated empirically 1000
143 tokens or words to be a reasonable size for learning the model of a document
144 and to be effective in its compression.

145 **3. Experimental Results**

146 The FCD as described in the previous section can be effectively employed
147 in tasks like authorship and plagiarism analysis. We report in this section
148 experiments on four datasets written in English, Italian, and German.

149 *3.1. The Federalist Papers*

150 We consider a dataset of English texts known as Federalist Papers, a col-
151 lection of 85 political articles written by Alexander Hamilton, James Madi-
152 son and John Jay, published in 1787-88 under the anonymous pseudonym
153 ‘Publius’. This corpus is particularly interesting, as Hamilton and Madison
154 claimed later the authorship of their texts, but a number of essays (the ones
155 numbered 49-58 and 62-63) have been claimed by both of them. This is a
156 classical dataset employed in the early days of authorship attribution liter-
157 ature, as the candidate authors are well-defined and the texts are uniform

$D(x)$	with equal	within the	word the	yet there
	with greater	within the union	worse than	yet what
	with many	within their	would at	you make
	with mutual	within which	would be	you render
[...]	with personal	without property	would be differently	you render him
	with success	without taking	would be unwise	you take
	with which	without violating	would certainly	zeal for
	with which they	without which	yet the	zeal in

Figure 1: Subset from a dictionary $D(x)$ extracted from a sample text x belonging to the Federalist dataset.

158 in thematics (Stamatatos, 2009). Several studies agreed on assigning the
159 disputed works in their entirety to Madison, while Papers 18-20 have gener-
160 ally been found to be written jointly by Hamilton and Madison as Hamilton
161 claimed, even though some researchers tend to attribute them to Madison
162 alone (Jockers and Witten, 2010; Meyerson, 2008; Adair, 1974).

163 We analyzed a dataset composed of a randomly selected number of texts
164 of certain attribution by Hamilton and Madison, plus all the disputed and
165 jointly written essays. We then computed a distance matrix related to the
166 described dataset according to the FCD distance, and performed on the
167 matrix a hierarchical clustering which is by definition unsupervised. A den-
168 drogram (binary tree) is heuristically derived to represent the distance ma-
169 trix in 2 dimensions through the application of genetic algorithms (Cilibrasi,
170 2007; Cilibrasi and Vitányi, 2005). Results are reported in Fig. 2, and
171 have been obtained using the freely available tool CompLearn available at
172 Cilibrasi et al. (2002). Each leaf represents a text, with the documents which
173 behave more similarly in terms of distances from all the others appearing as
174 siblings. The evaluation is done by visually inspecting if texts written by the
175 same authors are correctly clustered in some branch of the tree, i.e. by check-
176 ing how well the texts by the two authors can be isolated by ‘cutting’ the tree
177 at a convenient point. The clustering agrees with the general interpretation of
178 the texts: all the disputed texts are clearly placed in the section of the tree
179 containing Madison’s works. Furthermore, the three jointly written works
180 are clustered together and placed exactly between Hamilton and Madison’s

181 essays. We compare results with the hierarchical clustering derived from the
182 distance matrix obtained on the basis of NCD distances (Fig. 3), run with
183 the default blocksort compression algorithm provided by CompLearn: in this
184 case the misplacements of the documents is evident, as disputed works are
185 in general closer to Madison texts but are scattered throughout the tree.

186 3.2. *The Liber Liber dataset*

187 The rise of interest in compression-based methods is in part due to the
188 concept of relative entropy as described in Benedetto et al. (2002), which
189 quantifies a distance between two isolated strings relying on information the-
190 oretical notions. In this work the authors successfully perform clustering and
191 classification of documents: one of the considered problems is to automati-
192 cally recognize the authors of a collection comprising 90 texts of 11 known
193 Italian authors spanning the centuries XIII-XX, available at Onlus (2003).
194 Each text x was used as a query against the rest of the database, its clos-
195 est object y minimizing the relative entropy $D(x, y)$ was retrieved, and x
196 was then assigned to the author of y . In the following experiment the same
197 procedure as Benedetto et al. (2002) and a dataset as close as possible have
198 been adopted, with each text x assigned to the author of the text y which
199 minimizes $FCD(x, y)$. We compare our results with the ones obtained by
200 the Common N-grams (CNG) method proposed by Kešelj et al. (2003) us-
201 ing the most relevant 500, 1000 and 1500 3-grams in Table 2. The FCD
202 finds the correct author in 97.8% of the cases, while the best n-grams setting

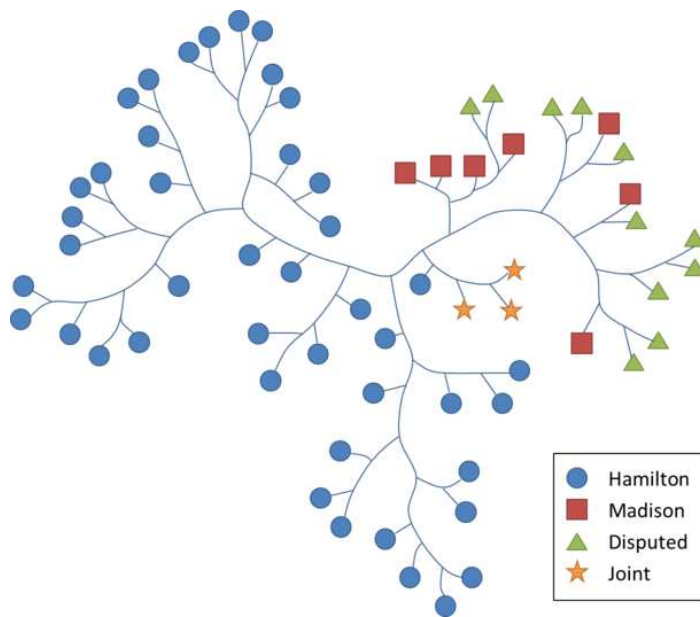


Figure 2: Hierarchical clustering of the Federalist dataset, derived by a full distance matrix obtained on the basis of the FCD distance.

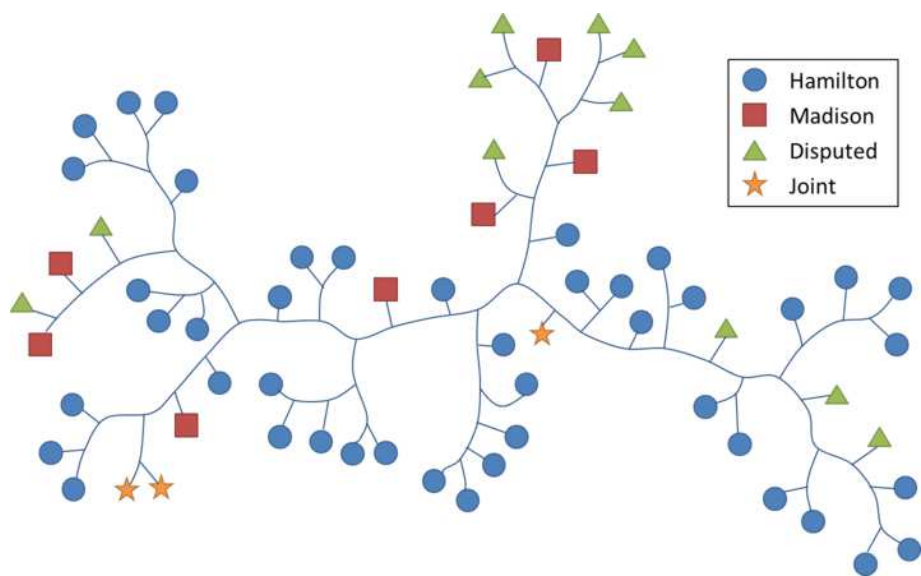


Figure 3: Hierarchical clustering of the Federalist dataset obtained on the basis of the NCD distance.

203 yields an accuracy of 90%. For FCD only two texts, *L'Asino* and *Discorsi*
204 *sopra la prima deca di Tito Livio*, both by Niccoló Machiavelli, are incor-
205 rectly assigned respectively to Dante and Guicciardini, but these errors may
206 be justified: the former is a poem strongly influenced by Dante (Caesar,
207 1989), while the latter was found similar to a collection of critical notes on
208 the very *Discorsi* compiled by Guicciardini, who was Machiavelli's friend
209 (Machiavelli et al., 2002). The N-grams-based method also assigns incor-
210 rectly Guicciardini's notes and a Dante's poem to Machiavelli, among others
211 misclassifications.

212 We also compared our results with an array of other compression-based
213 similarity measures (Table 3): our results outperform both the Ziv-Merhav
214 distance (Pereira Coutinho and Figueiredo, 2005) and the relative entropy as
215 described in Benedetto et al. (2002), while the algorithmic Kullback-Leibler
216 divergence (Cerra and Datcu, 2011) obtains the same results in a consider-
217 ably higher running time. Accuracy for the NCD method using an array
218 of linear compressors ranged from the 93.3% obtained using the bzip2 com-
219 pressor to the 96.6% obtained with the blocksort compressor. Even though
220 accuracies are comparable and the dataset may be small to be statistically
221 meaningful, another advantage of FCD over NCD is the decrease in compu-
222 tational complexity. While for NCD it took 202 seconds to build a distance
223 matrix for the 90 pre-formatted texts using the zlib compressor (with no
224 appreciable variation when using other compressors), just 35 seconds were
225 needed on the same machine for the FCD: 10 to extract the dictionaries and

Table 2: Classification results on the Liber Liber dataset. Each text from the 11 authors is used to query the database, and it is considered to be written by the author of the most similar retrieved work. The authors’ full names: Dante Alighieri, Gabriele D’Annunzio, Grazia Deledda, Antonio Fogazzaro, Francesco Guicciardini, Niccoló Machiavelli, Alessandro Manzoni, Luigi Pirandello, Emilio Salgari, Italo Svevo, Giovanni Verga. The CNG method has been tested using the reported amounts of n-grams.

Author	Texts	FCD	CNG-500	CNG-1000	CNG-1500
Dante Alighieri	8	8	6	5	7
D’Annunzio	4	4	4	3	4
Deledda	15	15	15	15	14
Fogazzaro	5	5	4	5	5
Guicciardini	6	6	5	5	5
Machiavelli	12	10	8	10	9
Manzoni	4	4	4	4	4
Pirandello	11	11	5	10	8
Salgari	11	11	10	10	9
Svevo	5	5	4	5	5
Verga	9	9	6	9	8
Total	90	88	71	81	78
Accuracy (%)	100	97.8	78.9	90	86.7

226 the rest to build the full distance matrix.

227 *3.3. The PAN Benchmark Dataset*

228 We tested our algorithm on datasets from the two most recent PAN CLEF
 229 (2013) competitions, which provide benchmark datasets for authorship attri-
 230 bution. From PAN 2013 we selected the author identification task described
 231 in Juola and Stamatatos (2013). In this task 349 training texts are provided,
 232 divided in 85 problems out of which 30 are in English, 30 in Greek and 25
 233 in Spanish. For each set of documents written by a single author it must be
 234 determined if a questioned document was written by the same author or not.
 235 Each text is approximately 1000 words long, which is close to our empirical

Table 3: Accuracy and running time for different compression-based methods applied to the Liber Liber dataset.

Method	Accuracy (%)	Running Time (sec)
FCD	97.8	35
Relative Entropy	95.4	NA
Ziv-Merhav	95.4	NA
NCD (zlib)	94.4	202
NCD (bzip2)	93.3	198
NCD (blocksort)	96.7	208
Algorithmic KL	97.8	450

236 estimation of the minimum size for FCD to find relevant patterns in a data
 237 instance (Section 2). For each problem, we consider an unknown text to be
 238 written by the same author of a given set of documents if the average FCD
 239 distance to the latter is smaller than the mean distance from all documents
 240 of a given language. Compared to the performance of the 18 methods re-
 241 ported in Juola and Stamatatos (2013), the FCD finds the correct solution
 242 in 72.9% of the cases and yields the second best results, ranking first for
 243 the set of English problems and fifth for both the Greek and Spanish sets
 244 (Table 4), outperforming among others two compression-based and several
 245 n-grams-based methods. It must be stressed that the FCD took approxi-
 246 mately 38 seconds to process the whole dataset, while the imposters method
 247 by Seidman (2013), which ranked first in the competition for all problems
 248 excluded the ones in Spanish, took more than 18 hours. Furthermore, the
 249 latter method requires the setting of a threshold, while the FCD skips this
 250 step. On the other hand, the contest participants had only a small subset of
 251 the available ground truth to test their algorithms.

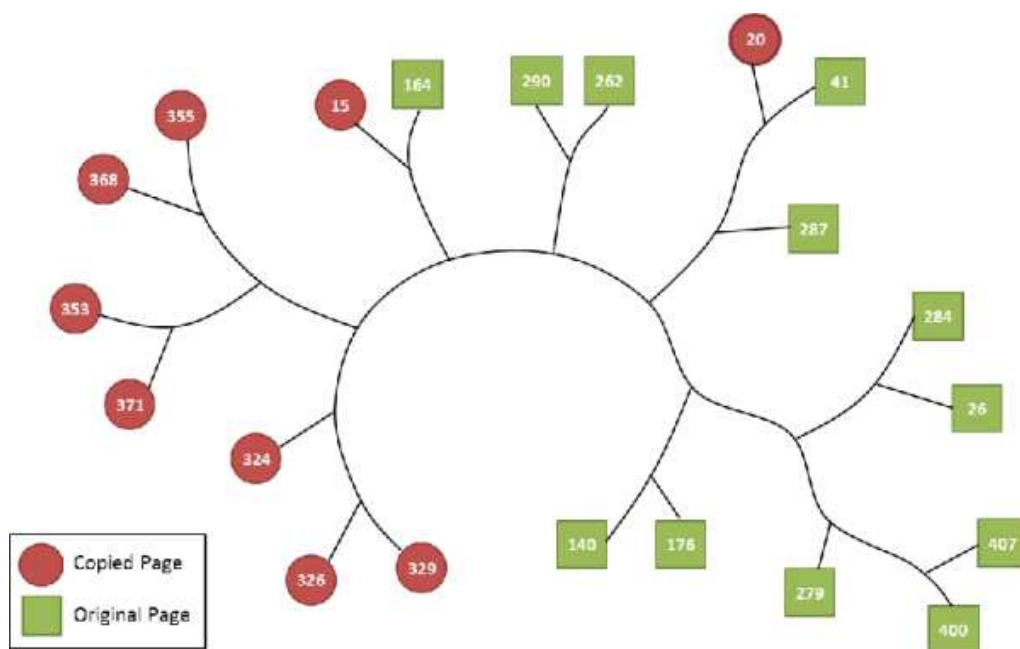


Figure 4: Hierarchical clustering of pages extracted from Guttenberg PhD thesis.

Table 4: Author identification task of the CLEF PAN 2013 dataset. The dataset contains 349 training texts plus 85 test documents of questioned authorship, with problems given in English, Greek and Spanish. The table reports how the FCD ranks compared to 18 participants to the PAN 2013 contest. The first ranked submission for each problem is reported as ‘Best PAN’.

Task	FCD	Best PAN	Rank
Overall	72.9 %	75.3 %	2
English	83 %	80 %	1
Greek	63 %	83 %	5
Spanish	72 %	84 %	5

252 We tested FCD also on the largest closed-class classification problem
 253 (task I) from the 2012 PAN competition: open-class problems were not
 254 considered as the simple classification algorithm adopted does not allow a
 255 rejection class. Using a corpus of 14 test and 28 training texts belonging
 256 to 14 different authors, the FCD (using a simple nearest neighbour classi-
 257 fication criterion) assigns correctly 12 out of 14 documents to their correct
 258 authors. Out of the 25 which took part to the competition, only 4 methods
 259 submitted by three groups (Sapkota and Solorio, 2012; Tanguy et al., 2012;
 260 Popescu and Grozea, 2012) outperformed our method (all of them with 13
 261 documents correctly recognized). As a comparison, the NCD and trigrams-
 262 based CNG (using the most meaningful 1000 trigrams per document, as this
 263 setting yields the best results in Table 2) assigned 2 and 9 documents out
 264 of 14 to the correct author, respectively. The results in Tables 4 and 5 are
 265 encouraging, specially if we consider that the FCD is a general method which
 266 is not specific for the described tasks.

Table 5: Classification results on task I of the CLEF PAN 2012 dataset. The dataset contains 28 texts belonging to 14 different authors for training and 14 for testing. The best results obtained in the PAN 2012 contest are reported as ‘Best’.

Method	FCD	NCD	CNG	Best
Correct (out of 14)	12	2	9	13

267 *3.4. The Guttenberg Case*

268 In February 2011, evidence was made public that the former German
 269 minister Karl-Theodor zu Guttenberg had violated the academic code by
 270 copying several passages of his PhD thesis without properly referencing them.
 271 This eventually led to Guttenberg losing his PhD title, resigning from being
 272 minister, and being nicknamed Baron Cut-and-Paste, Zu Copyberg and Zu
 273 Googleberg by the German media (BBC, 2011). Evidence of the plagiarism
 274 and a detailed list of the copied sections and of the different sources used by
 275 the minister is available at GuttenPlag (2011).

276 We selected randomly two sets of pages from this controversial disserta-
 277 tion, with the first containing plagiarism instances, and the second material
 278 originally written by the ex-minister. Then we performed an unsupervised
 279 hierarchical clustering on the distance matrix derived from FCD distances as
 280 described in Section 3.1. First attempts made by analyzing single pages failed
 281 at separating the original pages in a satisfactory way, as the compressor needs
 282 a reasonable amount of data to be able to correctly identify shared patterns
 283 between the texts. We selected then two-pages long excerpts from the thesis,
 284 with the resulting clustering reported in Fig. 4 showing a good separation of
 285 the texts containing plagiarism instances (in red in the picture). The only

286 confusion comes from pages starting at 41 with pages starting at 20, in the
287 bottom-left part of the clustering. This is justified by the fact that page 41
288 refers to the works of Loewenstein, who happens to be the same author from
289 which part of page 20 was plagiarized (Loewenstein, 1959). Therefore, the
290 system considers page 20 to be similar to the original style of the author at
291 page 41.

292 Even though the described procedure is not able to detect plagiarism, it
293 can find excerpts in a text which are similar to a given one. If instances of
294 plagiarized text can be identified, objects close to them in the hierarchical
295 clustering will be characterized by a similar style: therefore, this tool could
296 be helpful in identifying texts which are most likely to have been copied from
297 similar sources.

298 **4. Conclusions**

299 This paper evaluates the performance of compression-based similarity
300 measures on authorship and plagiarism analysis on natural texts. Instead of
301 the well-known Normalized Compression Distance (NCD), we propose using
302 the dictionary-based Fast Compression Distance (FCD), which decomposes
303 the texts in sets of reoccurring combinations of words captured in a dictio-
304 nary, which describe the text regularities, and are compared to estimate the
305 shared information between any two documents. The reported experiments
306 show the universality and adaptability of these methods, which can be ap-
307 plied without altering the general workflow to documents written in English,

308 Italian, Greek, Spanish and German. The main advantage of the FCD with
309 respect to traditional compression-based methods, apart from the reduced
310 computational complexity, is that it yields more accurate results. We can
311 justify this with two remarks: firstly, the FCD should be more robust since it
312 performs a word-based analysis, focusing exclusively on meaningful patterns
313 which better capture the information contained in the documents; secondly,
314 the use of a full dictionary allows discarding any limitation that real compres-
315 sors have concerning the size of buffers and lookup tables employed, being
316 the size of the dictionaries bounded only by the number of relevant patterns
317 contained in the objects. At the same time, the data-driven approach typi-
318 cal of NCD is maintained. This allows keeping an objective, parameter-free
319 workflow for all the problems considered in the applications section, in which
320 promising results are presented on collections of texts in Italian, English, and
321 German.

322 **References**

- 323 Adair, D., 1974. *Fame and the Founding Fathers*. Indianapolis: Liberty Fund.
- 324 BBC, 2011. Germany's Guttenberg 'deliberately' plagiarised.
325 URL <http://www.bbc.co.uk/news/world-europe-13310042>
- 326 Benedetto, D., Caglioti, E., Loreto, V., 2002. Language trees and zipping.
327 *Physical Review Letters* 88 (4), 48702.
- 328 Caesar, M., 1989. *Dante, the critical heritage, 1314-1870*. Routledge.

- 329 Cerra, D., Datcu, M., 2011. Algorithmic relative complexity. *Entropy* 13 (4),
330 902–914.
- 331 Cerra, D., Datcu, M., 2012. A fast compression-based similarity measure with
332 applications to content-based image retrieval. *Journal of Visual Commu-
333 nication and Image Representation* 23 (2), 293 – 302.
- 334 Chen, X., Francia, B., Li, M., McKinnon, B., Seker, A., 2004. Shared infor-
335 mation and program plagiarism detection. *IEEE Transactions on Informa-
336 tion Theory* 50 (7), 1545–1551.
- 337 Cilibrasi, R., 2007. *Statistical inference through data compression*. Lulu.com
338 Press.
- 339 Cilibrasi, R., Cruz, A., de Rooij, S., Keijzer, M., 2002. *CompLearn*.
340 URL <http://www.complearn.org>
- 341 Cilibrasi, R., Vitányi, P. M. B., 2005. Clustering by compression. *IEEE
342 Transactions on Information Theory* 51 (4), 1523–1545.
- 343 CLEF, 2013. *PAN Lab 2012 & 2013: Uncovering Plagiarism, Authorship,
344 and Social Software Misuse*.
345 URL <http://pan.webis.de>
- 346 de Graaff, R., 2012. *Authorship attribution using compression distances*.
347 Master thesis, Leiden University.
348 URL <http://www.liacs.nl/assets/Bachelorscripties/2012-18RamondeGraaff.pdf>

- 349 Gordo, A., Perronnin, F., Valveny, E., 2013. Large-scale document image
350 retrieval and classification with runlength histograms and binary embed-
351 dings. *Pattern Recognition* 46 (7), 1898 – 1905.
352 URL <http://www.sciencedirect.com/science/article/pii/S0031320312005304>
- 353 Grieve, J., 2007. Quantitative authorship attribution: An evaluation of tech-
354 niques. *Literary and Linguistic Computing* 22 (3), 251–270.
355 URL <http://llc.oxfordjournals.org/cgi/doi/10.1093/llc/fqm020>
- 356 GuttenPlag, 2011. Collaborative documentation of plagiarism.
357 <http://de.guttenplag.wikia.com>.
- 358 Jockers, M. L., Witten, D. M., 2010. A comparative study of machine learning
359 methods for authorship attribution. *Literary and Linguistic Computing*
360 25 (2), 215–223.
- 361 Juola, P., Stamatatos, E., 2013. Overview of the Author Identification Task
362 at PAN 2013.
363 URL <http://www.clef-initiative.eu/documents/71612/3095ffc3-376b-40eb-af10-8251c>
- 364 Juola, P., 2006. Authorship attribution. *Found. Trends Inf. Retr.* 1 (3), 233–
365 334.
366 URL <http://dx.doi.org/10.1561/1500000005>
- 367 Keogh, E., Lonardi, S., Ratanamahatana, C., 2004. Towards parameter-free
368 data mining. In: *Proceedings of the tenth ACM SIGKDD international*
369 *conference on Knowledge discovery and data mining*. ACM, pp. 206–215.

- 370 Kešelj, V., Peng, F., Cercone, N., Thomas, C., 2003. N-gram-based author
371 profiles for authorship attribution. In: Proceedings of the Conference Pa-
372 cific Association for Computational Linguistics, PACLING. Vol. 3. pp.
373 255–264.
- 374 Koppel, M., Schler, J., Argamon, S., January 2009. Computational methods
375 in authorship attribution. *J. Am. Soc. Inf. Sci. Technol.* 60 (1), 9–26.
376 URL <http://dx.doi.org/10.1002/asi.v60:1>
- 377 Koppel, M., Schler, J., Argamon, S., 2011. Authorship attribution in the
378 wild. *Language Resources and Evaluation* 45, 83–94, 10.1007/s10579-009-
379 9111-2.
380 URL <http://dx.doi.org/10.1007/s10579-009-9111-2>
- 381 Loewenstein, K., 1959. Verfassungsrecht und verfassungspraxis der vere-
382 inigten staaten. *Enzyklopaedie der Rechts- und Staatswissenschaft.*
- 383 Macedonas, A., Besiris, D., Economou, G., Fotopoulos, S., 2008. Dictionary
384 based color image retrieval. *Journal of Visual Communication and Image*
385 *Representation* 19 (7), 464–470.
- 386 Machiavelli, N., Atkinson, J., Sices, D., 2002. *The Sweetness of Power:*
387 *Machiavelli’s Discourses & Guicciardini’s Considerations.* Northern Illinois
388 University Press.
- 389 Meyerson, M., 2008. *Liberty’s blueprint: how Madison and Hamilton wrote*

390 The Federalist Papers, defined the Constitution, and made democracy safe
391 for the world. Basic Books.

392 Oliveira, W., Justino, E., Oliveira, L. S., 2013. Comparing compression
393 models for authorship attribution. Forensic Science International 228 (1-
394 3), 100 – 104.

395 URL <http://www.sciencedirect.com/science/article/pii/S0379073813000923>

396 Onlus, L. L., 2003. the Liber Liber dataset. <http://www.liberliber.it>.

397 Pereira Coutinho, D., Figueiredo, M., 2005. Information Theoretic Text Clas-
398 sification using the Ziv-Merhav Method. Pattern Recognition and Image
399 Analysis 3523, 355–362.

400 Popescu, M., Grozea, C., 2012. Kernel methods and string kernels for au-
401 thorship analysis. In: Forner, P., Karlgren, J., Womser-Hacker, C. (Eds.),
402 CLEF (Online Working Notes/Labs/Workshop).

403 URL <http://dblp.uni-trier.de/db/conf/clef/clef2012w.html>

404 Sapkota, U., Solorio, T., 2012. Sub-profiling by linguistic dimensions to solve
405 the authorship attribution task. In: Forner, P., Karlgren, J., Womser-
406 Hacker, C. (Eds.), CLEF (Online Working Notes/Labs/Workshop).

407 URL <http://dblp.uni-trier.de/db/conf/clef/clef2012w.html>

408 Seidman, U., 2013. Authorship Verification Using the Impostors Method.
409 Notebook for PAN at CLEF 2013

410 URL <http://www.clef-initiative.eu/documents/71612/7a4e6a71-46e9-4bb1-ab66-8ea9c>

- 411 Stamatatos, E., 2009. A survey of modern authorship attribution methods.
412 J. Am. Soc. Inf. Sci. 60 (3), 538–556.
- 413 Tanguy, L., Sajous, F., Calderone, B., Hathout, N., 2012. Authorship at-
414 tribution: Using rich linguistic features when training data is scarce. In:
415 Forner, P., Karlgren, J., Womser-Hacker, C. (Eds.), CLEF (Online Work-
416 ing Notes/Labs/Workshop).
417 URL <http://dblp.uni-trier.de/db/conf/clef/clef2012w.html>
- 418 Watanabe, T., Sugawara, K., Sugihara, H., 2002. A new pattern representa-
419 tion scheme using data compression. IEEE Transactions on Pattern Anal-
420 ysis and Machine Intelligence 24 (5), 579–590.
- 421 Welch, T., 1984. Technique for high-performance data compression. IEEE
422 Computer 17 (6), 8–19.