

The Comparative Power of *Type/Token* and *Hapax legomena/Type* Ratios: A Corpus-based Study of Authorial Differentiation

Sundus Muhsin Ali

College of Arts, Baghdad University

E-mail: dr_sun_alubaidy@yahoo.com

Khalid Shakir Hussein

College of Education, Thi-Qar University

E-mail: khalidshakir74@gmail.com

Doi:10.7575/aiac.all.v.5n.3p.112

Received: 04/04/2014

URL: <http://dx.doi.org/10.7575/aiac.all.v.5n.3p.112>

Accepted: 14/05/2014

Abstract

This paper presents an attempt to verify the comparative power of two statistical features: Type/Token, and Hapax legomena/Token ratios (henceforth TTR and HTR). A corpus of ten novels is compiled. Then sixteen samples (each is 5,000 tokens in length) are taken randomly out of these novels as representative blocks. The researchers observe the way TTR and HTR behave in discriminating four novelists: Joyce, Woolf, Faulkner and Hemingway. When compared to the traditional statistical features (e.g. word length average, Sentence length average, etc.), TTR and HTR are by far more competent in comparing the distinctive quantitative behavior of each novelist. It turns out that TTR and HTR contribute more or less in creating a sort of statistical identity which can be used in giving a vivid comparison and discrimination of the four novelists involved in this paper. Nevertheless, HTR sounds more viable in achieving the discriminating task than TTR.

Keywords: Statistical linguistics, Corpus linguistics, Stylometrics

1. Introduction

Differentiating authors on statistical grounds has a long standing history. The earliest technique to give a statistical distinction to a particular writer was suggested by Mendenhall (1887, 1901). The first marker targeted by Mendenhall was *word-length*. Yule (1944) suggested *sentence length* to be a potential standard for a method of authorship discrimination, though he made himself clear later on about the indecisive reliability of this method. Fuchs (1952) (cited in Holmes, 1994:88) worked on *the syllable* as another sensitive stylistic marker. In 1966, Sommers (cited in Holmes, 1994: 89) resorted to distribution of parts of speech in a text as a more reliable authorial marker. The historical list of these statistical markers looks endless and lies beyond the scope of this paper.

Authors' styles can be assigned numerical measures representing the frequency counts of the markers mentioned above. However, the unprecedented progress in computer science gave access to more subtle statistical explorations of textual corpora especially after the tremendous increase in the quantities of online linguistic data. Thus, the computational progress together with the availability of machine-readable bodies of data gave birth to a new set of statistical measures that have not been checked before this revolution in technology. TTR and HTR are among the newly born measures that became accessible to researchers. The question raised by the researchers in this paper is: *how powerful are they in discerning differences among authorial styles?* The following section will overview the statistical basics of TTR and HTR.

1.1 TTR

TTR is basically one particular statistical strategy of *tokenization* that encompasses a given set of basic statistical measures used in looking at texts with computerized methods. Whatever was the strategy used in tokenizing a digital text, it always utilizes the same statistical concepts of *tokens* and *types*. A token is any single linguistic unit, most often a word, in a text (Baker et al, 2006: 159). While the number of tokens in a computerized database refers to the total number of words. As for the number of types, it refers to the total number of the unique distinct type of words (ibid: 162). Therefore, a token is any linguistic item that occurs in a text regardless of its type, whereas a type is a statistical concept that targets only the token-types involved in a surveyed corpus. Comparing the number of *tokens* in the data or corpus to the number of *types of tokens* can tell us much about how large a range of vocabulary is used in the corpus under consideration. For example, the sentence, (The man put the book on the table), contains (eight) tokens (*the, man, put, the, book, on, the, and table*) but only (six) types (*the, man, put, book, on, and table*). So its TTR is simply $(6 \div 8 = 0.75)$ types per token.

As a text gets longer, the number of types that have already been encountered increases, and the likelihood of any given token representing a new type goes down (Hardie & McEnery, 2006:139). For this reason, TTR as a statistic measure is very sensitive to the length of the text under investigation. Therefore, it is extremely crucial when comparing the TTRs of different texts to make sure that the texts concerned are of equal sizes (Baker et al, 2006:138). Generally speaking, a high TTR indicates a large amount of *lexical variation* and a low TTR indicates relatively little *lexical variation* (ibid.).

As with *lexical density* or *richness*, referring to a measure calculated by counting the number of lexical words in a text that occur only once, the TTR can also be used to monitor changes in the use of vocabulary items throughout one text or throughout a group of texts produced by one person (Butler, 1995: 135).

1.2 HTR

There is a whole series of criteria, linguistic resources and discourse strategies suggested by linguists to capture the authorial blueprint: *the degree of unity, completeness and coherence/cohesion* of the texts under comparison, *inconsistency in referential style, decontextualization, . . .* etc. (see Turell, 2008: 282-7). However, the researchers are particularly interested in one original marker: *hapax-legomena*. This does not mean that the other markers are irrelevant but they are rather reachable only through certain specialized packages of programs such as *Vocalyse* and *File Comparison*. These are two programs written by Woolls and then developed into *Copycatch Gold* (2002). Packages like these made it possible for linguists to measure up certain statistical markers seen as being most critical and direct in exploring the statistical behavior of a given corpus. Such markers range from *shared vocabulary*, *unique vocabulary*, to *shared once-only words* (ibid: 288).

This paper, nevertheless, relies solely on *WordSmith Tools*, a package designed basically to conduct statistical processing of texts and not to detect plagiarism or doubtful cases of statistical . Some outputs of WordSmith Tools can imply particular indications as a sort of byproduct pieces of information. *Wordlists*, for example, bring in, amongst other things, *hapax-legomena* as byproduct markers figured out in such word frequency counts.

The contribution of the notion of *hapax-legomena* to discriminate authorial identities is deeply entrenched in their low distribution or infrequent occurrence. A hapax is commonly defined as a word that occurs only *once* in a text or a corpus (Lardilleux & Lepage, 2007: 458). The distinctive power of hapaxes comes about from the parallelism between their unique occurrence, on the one hand, and the author's linguistic uniqueness or *idiolect*, on the other.

If two texts are available for comparison, the vocabulary used in both texts can be compared *quantitatively* so that a particular "qualitative textual analysis" can be backed up (Johnson, 1997), and the amount of the shared hapaxes is supposed to be quite powerful in establishing authorial distinctiveness: *the more there is, the less the authorial differentiation*.

It might be plausible to avoid restricting authorial differentiation to particular markers instead of surveying more varying features which might be lexical, grammatical and textual. This could contribute more effectively in figuring out how far the texts in question are similar or different enough to be attributed to one or different authors. However, the more markers involved, the more programs should be used, and not to mention the indispensable manual checking that some of these features might be in need. Besides, most of these features are case-dependent and rarely do they have stable and consistent confirmatory results all through the textual bodies they are scored in.

Things are rather different with hapaxes: a high number of *shared* hapax legomena and a low number of *unique* hapax legomena are both quite useful in measuring up how likely the authorial differentiation is (cited in Coulthard et al., 2010: 526). So, what is looked for here is the particular type of the items shared. Those items shared *once only* (hapaxes) are supposed to be *insignificant* in terms of the main concern of the text. This would explain why they have been used less frequently, actually just once. Otherwise, they would occur more frequently rather than showing up for once (Coulthard, 2004: 5). Thus the hypothesis is made so clear by Coulthard (ibid.) in this concern:

The chances of two writers independently choosing several of the same words for single use are so remote as to be discountable.

What is more, the researchers think that there is a relatively stable ratio of hapax that can be safely and distinctively attributed to authors throughout their writings. This ratio compares the number of hapaxes to the number of *types* (not *tokens*) along the text itself. It might be a quite useful marker in raising doubts about the occurrence of plagiarism as a case of authorial indifferenciation. Statistically speaking, any extreme variation in the figures alluding to the ratio of hapaxes to types is most likely to alert suspicion of plagiarism or doubtful authorial identity.

Despite the negative attitude most statisticians have to hapaxes since they used to neglect those tokens and types attaining low frequency, hapaxes in particular constitute something like (40%) of words used in a corpus (Lardilleux & Lepage, 2007: 458). Beside the type/token ratio, the (HTR) enhances the indication of vocabulary richness. The more hapaxes are there, the richer the vocabulary. This can give a measure of the author's distinctive originality that can bring him into a sort of statistical discrimination. Throughout this study it will hopefully be seen how far hapaxes can be taken as reliable markers of author's linguistic identity.

2. A Summary of the Analysis Procedures

The corpus used in this paper is compiled via the Internet. As for the digital samples included in the corpus, they will be subjected to *five* analysis procedures:

1. *Authenticity Investigation*: it is quite expected that "a corrupt sample" would most definitely produce "a corrupt analysis". Consequently, it should be determined that the digital samples selected for authorship analysis are clean samples, a task which sounds extremely difficult if not impossible (Juola, 2008:247). Since the samples selected for this study are machine-readable, the scanning or retyping processes could be a very threatening source of all types of errors. The researchers tried their best to check the authenticity of each sample making sure that each one is highly representative of the authors involved. Whenever there are hard copies of the texts they should be compared to their digital ones. This process might appear boring and painstaking but it is inescapable. Moreover, there are certain "non-authorial" materials that should be removed from the samples: *major heads, section heads, page numbers, quotations*, and so forth. They could be a severe threat to the statistics ascribed to the author's linguistic habits. Juola (ibid: 248) puts it quite directly:

. . . , all extraneous material that did not come from the author's pen (or keyboard) should be eliminated, a task requiring extreme care and knowledge on the part of the researcher . . .

Hence, only the main body of the texts will be considered: *titles, author names, dates, . . .*etc. all were excluded. After all, every sample should be authenticated in a way that sounds independent and reliable. Otherwise, the sample would be eliminated for its potential extraneous variables that might influence the statistical results

2. Transcribing samples into *plain text format*

3. Grouping all the samples into one master corpus

4. Analyzing samples with their master corpus via *WordSmith Tools (5.0)* for frequency and word count, besides producing some sort of charts representing basic statistical descriptions

5. Importing WordSmith Tools *outputs* into an excel spreadsheet in a form of matrix.

3. An Experiment in Authorial Differentiation

In order to evaluate the performance of the TTR and HTR methods concerned in this paper the researchers perform a classification experiment. It sounds crucial to check out the discriminating power of a statistical methodology before conducting it in real-life cases of authorial differentiation or characterization.

The experiment has been conducted on (10) novels: *Ulysses*, and *Finnegans Wake* by James Joyce; *Mrs Dalloway*, *To the Lighthouse*, and *The Waves* by Virginia Woolf; *As I Lay Dying*, *Light in August*, and *The Sound and the Fury* by William Faulkner; and lastly *A Farewell To Arms*, and *The sun Also Rises* by Ernest Hemingway.

Sixteen samples were selected randomly from the novels above. Table (1) below shows a full description of the corpus.

Table 1. Distribution of English Corpus by Genre

Author	Text-samples	Samples Number	Genre
James Joyce	(5,000-token) samples selected from two novels <i>Ulysses</i> and <i>Finnegans Wake</i> .	4	Fiction
Virginia Woolf	(5,000-token) samples selected from three novels <i>Mrs Dalloway</i> , <i>To the Lighthouse</i> , and <i>The Waves</i> .	4	Fiction
William Faulkner	(5,000-token) samples selected from three novels <i>As I Lay Dying</i> , <i>Light in August</i> , and <i>The Sound and the Fury</i> .	4	Fiction
Ernest Hemingway	(5,000-token) samples selected from two novels <i>A Farewell To Arms</i> , and <i>The Sun Also Rises</i> .	4	Fiction

The four novelists, however, haven not been selected haphazardly but rather purposively. Joyce' writings in particular represent a very threatening challenge to the statistical method adopted in this study. Joyce had unprecedented desire to experiment with language and this desire shows itself most intensively throughout his two controversial novels: *Ulysses* and *Finnegans Wake* (see Tadie, 2003). Investigating the various manifestations of Joyce's linguistic experimentations is by far beyond the scope of this study. However, what sets a quite challenging threat to the validity of the method used in this study is Joyce's deviations from the conventional morphological patterns. This aspect of linguistic experimentation could hamper the mechanism used in WordSmith Tools that relies on processing words into distinct *tokens* and *types* to produce clear cut wordlists.

Joyce's usage of words, in the morphological sense, defies all the preconceived ideas about words as being the counterpart of bricks in a linguistic construction (see Cordell,1997). The morphological techniques used by him have produced a quite bizarre set of words or tokens, so to speak: "Thursdaymomun" (*Ulysses*,2013: 536); "Nationalgymnasiummuseumsanatoriumandsuspensoriumsordinaryprivatedocent-generalhistoryspecialprofessor" (ibid:292),"Lukkedoerendunandurraskewdylooshoofermoypuertooryzoosphalnabortsportthaokansakroidverjkapakk

apuk" (Finnegans Wake, 2013: 27-28). . .etc. He even used to condense dozens of function words in one queer and hybrid construction deforming the graphological independence of these words. This kind of radical morphological innovations would most definitely raise a question that should be urgently answered: *have Joyce's morphological experimentations influenced his authorial consistency? Has this kind of bizarre tokens and types been used in a statistically consistent way that would not affect the statistical stability of the TTR and HTR in the discriminating process? Or they have constituted a sort of statistical turbulence that hampered the viability of the discriminating method?*

As for Woolf and Faulkner, they set another challenge that the method need to address. It is a well-established point that they both have their own distinct styles, though they belong to the same technical school of stream of consciousness. Their styles are said to be overlapping in more than one area of linguistic usage. The question posed in this context is *how far does this so-called stylistic overlapping have repercussions on the discriminating power of the TTR and HTR?* This question stirs some predictions related to two possibilities : their authorial consistencies might be either so distinct so that they can be attributionally discriminated or extremely and unattributionally look alike a matter which could destroy the concept of authorial originality.

Hemingway, on the other hand, does not belong to the technical school of stream of consciousness and hence does not share any of the linguistic traits usually ascribed to Joyce, Faulkner, and Woolf. He has his own well-known stylistic traits and markers that leave the reader with an impression of reading a sort of simple and journalistic language. *Does that give him an advantage over the other novelists of having an authorial consistency that is radically distinct from that claimed for the novelists in question?*

These three challenges or questions have been utterly addressed in this experiment which is basically conducted to check the consistency of the writing style of each novelist. What is generally questioned here is *whether it is valid to hypothesize that a novelist maintains his/her distinctive authorial identity regardless of how experimental his/her language is or the literary school he/she belongs to.*

4. Results (Basic Statistics)

The results of the statistical processing of the master corpus being compared to the statistics of the samples that constitute its whole body (subcorpora) are shown along the five analysis procedures summarized above.

After conducting a thorough authentication of the samples and comparing them to their equivalent hard copies, the researcher modified each corpus by removing common contractions, especially those involving function words such as: can't – can not; won't – will not; couldn't – could not; you're – you are . . .etc. To determine the hierarchical list of the function words the researcher combined the sixteen samples into one master corpus. Table (5.1) shows the basic statistical descriptions of the combined corpus so that one can check out the discriminating power of the traditional statistical features.

It is quite obvious that word-length and sentence-length are utterly indiscriminative and can not be used as potential characteristics of authorship. Though these two features have been used in the past as being reliable markers of the writing style, which might be true to some extent, however, they lose their decisive weight in the experiment under consideration. Table (1) below provides us with an evidence that they do not reflect authorial distinction of any type and can not be used as trustworthy determiners of the authorial identity. The mean of sentence-length in words ranges from (7.71) up to (23.57). This is a rather wide range and no novelist has shown any stable statistics for his sentences length.

The same can be concluded on word-length but with an exceptional difference that the range is utterly narrow: it fluctuates between (3.88) and (4.84). Nevertheless, this relatively stable statistics can not be exclusively ascribed to any particular author. It might be reasonably possible to assume this statistics as a trait of English in general. This could explain the reason why the samples show the harmonious yet indiscriminative statistics.

Unlike the statistical features in the table below, TTR holds an exceptional position since it allows us to conclude little about authorship and authorial differentiation. There are crucial distinctions between the TTRs calculated for every novelist indicating rather promising guesses about the lexical diversity (or vocabulary richness) of each author. Joyce interestingly is by far the richest in his vocabulary scoring a TTR value rising up to (51.9), while Faulkner and Woolf, though distinctive, come near each other with a minimum TTR (21.6) and maximum value (26.3). Still, Woolf's vocabulary appears to be distinctively richer than Faulkner's. As shown in Figure (1), Hemingway's TTR comes down holding comparatively poor lexical diversity. With a style said to lack substance with little description of emotion, it is extremely expected that Hemingway's vocabulary is limited and can be readily distinguished from that of the other novelists. Therefore, it is attributed to a statistical area thrown distinctively away from that occupied by Faulkner and Woolf.

The best that one can conclude about the numerical values of TTR in Table (2) below, in so far as authorial discrimination is concerned, is that Joyce's samples stand out in a quite privileged way. That is why their authorship can be safely attributed to Joyce as a result of their own extraordinary lexical diversity. Whereas Woolf, Faulkner, and Hemingway's TTRs are not trustworthy as they attain a low set of figures suggesting a lot of repeated lexical items (or tokens) all through the samples that belong to them.

Table 2. Basic Statistics of the Master Corpus

Authors statistics	Joyce				Faulkner				Woolf				Hemingway			
	Ulysses		Finnegans		As I Lay Dying		Lgt. in Aug.	The S F	Mrs Dalloway		The Waves	To Light H	A Farewell	A Farewell	The Sun Also	The Sun Also
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
File Size	29,857	28,89	29,391	30,88	26,155	25,829	27,390	26,270	27,617	28,227	27153	26649	25900	26,652	26,296	26,907
tokens	5,104	5,020	5,138	5,123	5,030	5,016	5,090	5,021	5,079	5,063	5,051	4,919	5,140	5,178	5,052	5,011
types	2,340	1,988	2,483	2,662	1,139	1,139	1,103	1,238	1,153	1,319	1,262	1,294	1,028	1,018	1,058	1,010
type/token ratio (TTR)	45.89	39.60	48.33	51.97	22.71	22.71	21.67	24.66	22.70	26.05	24.99	26.31	20.02	19.66	20.95	20.16
standardised TTR	58.52	51.50	58.31	60.70	38.06	38.80	37.52	40.66	38.42	41.88	40.09	41.15	35.24	34.34	36.20	36.88
standardised TTR std.dev.	35.73	40.32	34.17	31.73	50.84	49.69	50.49	49.95	49.34	47.12	49.08	47.12	53.11	54.32	51.98	51.18
mean word length (in characters)	4.54	4.52	4.51	4.84	3.91	3.88	4.07	4.05	4.18	4.35	4.19	4.23	3.83	3.91	3.89	3.99
word length std.dev.	2.51	2.58	2.60	2.87	1.85	1.87	2.04	1.96	2.07	2.19	2.09	2.14	1.94	2.05	1.88	1.95
sentences	570	213	281	363	392	395	384	413	397	311	307	238	509	558	576	650
mean (in words)	8.95	23.57	18.28	14.11	12.80	12.70	13.26	12.16	12.79	16.28	16.44	20.66	10.09	9.28	8.77	7.71

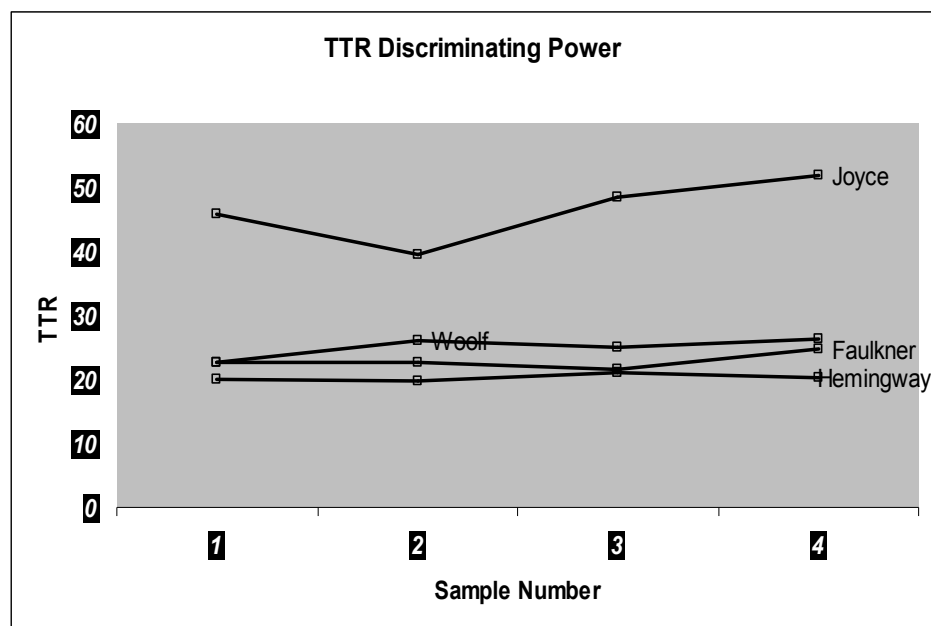


Figure 1. TTR Discriminating Power

Nevertheless, the writers' authorial identities acquire rather more discriminative power when measuring the values of HTR (see Table 3). The difference one can see in Figures (1) and (2) might be insignificant in Joyce's case since he is still distinguished with an exceptionally genuine authorial identity. HTRs have, however, brought about unavoidable differences that help much in making more reasonable predictions about Woolf, Faulkner, and Hemingway's authorial discriminations.

Table 3. HTR calculated for the Master Corpus

Author	Segment	Types	Hapax Legomena	HTR	Average
Joyce	1	2661	2391	0.89	0.84
	2	2482	2189	0.88	
	3	1988	1580	0.79	
	4	2340	1886	0.80	
Woolf	1	1153	682	0.59	0.62
	2	1319	829	0.62	
	3	1293	842	0.65	
	4	1262	790	0.62	
Faulkner	1	1139	632	0.55	0.56
	2	1103	599	0.54	
	3	1238	727	0.58	
	4	1196	649	0.54	
Hemingway	1	1028	535	0.52	0.51
	2	1018	523	0.51	
	3	1058	540	0.51	
	4	1010	528	0.52	

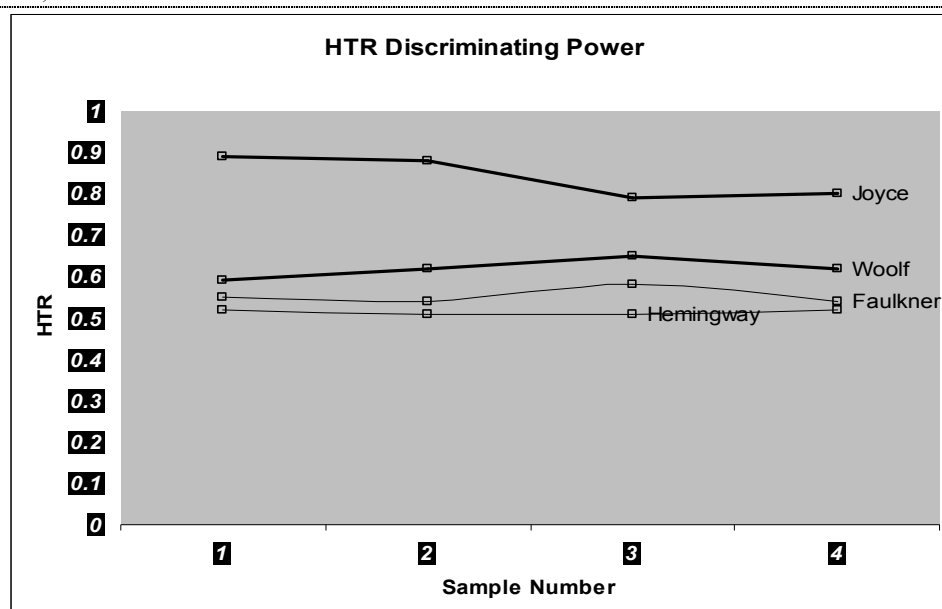


Figure 2.HTR Discriminating Power

5. Conclusions

The statistical analysis conducted throughout this paper has foregrounded the discriminative potentialities of TTR and HTR against all the traditional statistical measures. It turns out that HTR and TTR have generally good power for discriminating the writing styles of a difficult set of authors.

Even with Joyce's highly experimental language, they proved robust and stable. Joyce's writing style has always achieved a considerable degree of distinction when compared to others. The so-called case of stylistic overlapping between Woolf and Faulkner does not have serious consequences on the discriminating power of these two ratios. Though Woolf and Faulkner cluster close to each other, they are still statistically recognized as two authors each has his/her distinct authorial style. They might look quite similar but not indistinguishable. Hemingway's simple telegraphic language has quite consistent values of TTR and HTR so that he was statistically spotted in another discriminative area within the statistical environment.

However, the statistical behavior of HTR might look more sophisticated than that of TTR. A quizzical look at the quantitative representations of the samples would not escape the informative performance of HTR. Some minute and significant statistical differences might be overlooked when observing only the way TTR measures up the statistical consistencies of different authorial styles.

References

- Baker, P., Hardie, A., & McEnery, T. (2006). *A Glossary of Corpus Linguistics*. Edinburgh: Edinburgh University Press.
- Butler, C. (1995). *Statistics in Linguistics*. Oxford: Blackwell.
- Cordell, D. K. (1997). *The Word According to James Joyce*. Massachusetts: Associated University Press, Inc.
- Coulthard, M., & Johnson, A., Kredens, K., & Woolls, D. (2010). Plagiarism: four forensic linguists responses to suspected plagiarism. In Coulthard, M., and Johnson, A., *The Routledge Handbook of Forensic Linguistics*. London and New York: Routledge.
- _____ (2004). Author identification, idiolect and linguistic uniqueness. *Applied Linguistics*, 25(4), 431-47.
- Faulkner, W. *The Sound and the Fury*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (13 July 2013).
- _____ *As I lay Dying*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (17 July 2013).
- _____ (1972) . *Light in August*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (17 July 2013).

- Hardie, A. & McEnery, T. (2006). statistics. In Sarah, G. & William, J. (eds.). *Encyclopedia of Language and Linguistics*. Retrieved from <http://www.sciencedirect.com> (02 May 2013).
- Hemingway, E. *A Farewell To Arms*. Retrieved from (http://www.ernest_hemingway_ecel_zengi-eng.com) (22 July 2013).
- _____. *The Sun Also Rises*. Retrieved from (http://www.ernest_hemingway_ecel_zengi-eng.com) (22 July 2013).
- Johnson, A., (1997). Textual kidnapping: a case of plagiarism among three student texts. *International Journal of Speech, Language and the Law*, 4(ii), 210-25.
- Joyce, J. *Ulysses*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (01 July 2013).
- _____. *Finnegans Wake*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (01 July 2013).
- Juola, P., (2008). *Authorship Attribution: Foundations and Trends*. Boston-Delft: now Publishers.
- Lardilleux, A., & Lepage, Y., (2007). The contribution of the notion of hapax legomena to word alignment. In *Proceedings of the 4th Language and Technology Conference (LTC'07)*, 458-462.
- Tadie', B. (2003). *Joyce and Contemporary Linguistic Theories*. Cambridge: Cambridge University Press.
- Turell, M. (2008). Plagiarism. In Gibbons, J., and Turell, M. (eds.) *Dimensions of Forensic Linguistics*. Philadelphia: John Benjamins Publishing Company.
- Woolls, D., (2006). Plagiarism. In Brown, K. (ed.). *The Encyclopedia of Language and Linguistics*, 2nd ed, Vol. 9, Oxford: Elsevier.
- Woolf, V. (1977). *Mrs Dallaway*. Retrieved from (<http://ebooks.adelaide.edu.au/>) (10 July 2013).
- _____. (1972). *To the Lighthouse*. Retrieved from (<http://ebooks.adelaide.edu.au/w/woolf/virginia/w91t/part3.html>) (13 July 2013).
- _____. (1985). *The Waves*. Retrieved from (<http://ebooks.adelaide.edu.au/w/woolf/virginia/w91t/part3.html>) (09 July 2013).