

## **The relation between word length and sentence length: an intra-systemic perspective in the core data structure**

*Peter Grzybek<sup>1</sup>, Emmerich Kelih<sup>1</sup>, Ernst Stadlober<sup>2</sup>*

**Abstract.** Word length and sentence length are systematically organized in texts and corpora. In recent attempts at the synergetic modeling of the relation between sentence length and word length, the importance of distinguishing intra-textual from inter-textual approaches has been emphasized. The present study focuses on the intra-textual level: with a particular emphasis on different text types, it is shown, under which conditions processes of inter-level self-regulation are operative, and when they fail to be efficient.

*Keywords:* Menzerath-Altmann law, word length, sentence length, interrelation, intra-systemic structure

### **1 Theoretical ruminations**

The impact of word length (WL) and sentence length (SL) for purposes of text classification has been repeatedly documented (cf. Grzybek et al. 2005; Kelih et al. 2006; Antić et al. 2006). Extending these studies, Grzybek and Stadlober (2007) and Grzybek et al. (2007) have focused on the relationship between SL and WL, rather than on these two linguistic categories as separate phenomena in their own right.

In this context, the relevance of Arens' Law has been emphasized and submitted to some critical re-investigation. Arens' Law is an extension of the well-known Menzerath Law which, subsequent to its generalization and mathematical formulation by Altmann (1980) has also become known as Menzerath-Altmann Law. The latter aims at a theoretical description of the relation of linguistic units of different levels. Basically, it claims that the complexity or length of a particular (linguistic) component is a function of the length or complexity of the (linguistic) construct which it constitutes; it has been successfully applied in systems theoretical analyses other than linguistic as well (Altmann and Schwibbe 1989). The most general form of what is known today as the Menzerath-Altmann Law, has been suggested by Altmann (1980) in his seminal "Prolegomena to Menzerath's Law":

$$(1a) \quad y = ax^{-b}e^{cx} \quad (a, b, c > 0),$$

with two special cases for  $c = 0$ , or  $b = 0$ , respectively, namely

$$(1b) \quad y = ax^{-b}, \text{ and}$$

$$(1c) \quad y = ae^{cx}$$

Only recently, Wimmer and Altmann (2005, 2006) have extended this approach in their "General derivation of some linguistic laws". It is based on the differential equation

---

<sup>1</sup> Institut für Slawistik, Universität Graz, Merangasse 70, A-8010 Graz, Austria; correspondence address: [peter.grzybek@uni-graz.at](mailto:peter.grzybek@uni-graz.at)

<sup>2</sup> Institut für Statistik, Technische Universität Graz, Steyrergasse 17/IV, A-8010 Graz, Austria.

$$(2) \quad \frac{dy}{y} = \left( a_0 + \frac{a_1}{x} + \frac{a_2}{x^2} + \frac{a_3}{x^3} + \dots \right) dx$$

resulting in the solution

$$(3) \quad y = C e^{a_0 x} x^{a_1} e^{-a_2/x - a_3/(2x^2)} \dots$$

With  $a_3 = 0$  and  $-a_2 = d$  in equation (3), they arrive at the addition of an optional factor  $e^{d/x}$ , thus obtaining six options with  $d = 0$  for equations (1a-c), where  $C = a$ ,  $a_0 = c$ ,  $a_1 = -b$ :

$$(1d) \quad y = a e^{d/x}$$

$$(1e) \quad y = a x^{-b} e^{d/x}$$

$$(1f) \quad y = a x^{-b} e^{cx} e^{d/x}$$

Anyway, equation (1a) is generally considered the most basic and commonly used “standard form” for linguistic purposes. With  $b > 0$ , it predicts a **decrease** in length or complexity of the linguistic components with an increase in length or complexity of the construct they constitute – in longer words, e.g., the syllables forming these words are predicted to be shorter than those forming shorter words.

These ruminations are of course of central importance for the relation between sentence length and word length. However, in his analyses of German literary prose texts, Arens (1965) observed an **increase** in sentence length going along with an increase of word length, thus obtaining a result seemingly contradictory to the expectations.

By way of a solution, Altmann (1983: 31), in his attempt to interpret these results in Menzerathian terms, pointed out that the Menzerath-Altman Law as described above is likely to hold true only when one is concerned with the direct constituents of a given construct. In case of the SL-WL relation, however, an intermediate level may be assumed to come into play – such as, e.g., phrases or clauses as the direct constituents of the sentence. As a consequence, words might be seen as direct constituents of clauses or phrases, but only as indirect constituents of sentences. Therefore, in its direct form, the Menzerath-Altman Law might fail to grasp the SL-WL relation. In this case, an increase in SL should indeed result in an increase of WL, and it should be expected to be of the Menzerathian non-linear form: with  $y$  symbolizing word length,  $z$  symbolizing phrase (or clause) length, and  $x$  symbolizing sentence length, we were thus concerned with two simultaneous relations,  $y = a z^{-b} e^{cz}$  and  $z = a' x^{-b'} e^{c'x}$ . Inserting the latter equation into the first, one obtains  $y$  as a function

$$(4) \quad y = a'' x^{b''} \exp\left(-c''x + a''' x^{-b'} e^{c'x}\right)$$

However, in studies of direct relations between linguistic units of different levels, the “standard case” of the Menzerath-Altman Law, i.e.  $z = a' x^{-b'}$  and  $y = a z^{-b}$ , has been sufficient. Following this line, one thus obtains  $y = a'' x^{b''}$ , for the indirect relation between sentence length and word length, corresponding to equation (1b). From this perspective, Arens Law is a special case of the Menzerath-Altman Law: the only difference between direct and indirect relations thus is that, in case of directly neighboring units, the exponents  $-b$  and  $-b'$  are negative (due to the predicted decline), whereas in case of indirectly related units, with intermediate levels,  $b'' = (-b) \cdot (-b')$  will become positive. However, this would hold true only in case of deterministic relations, and in no case for averages.

## 2. Empirical findings

Despite the importance of Arens' Law for linguistic and non-linguistic analyses in the field of general systems theory, only few studies have explicitly referred to it. A possible reason for this might be that there seems to be only poor evidence in support of the theoretical assumptions, as recently pointed out by Grzybek and Stadlober (2007). Thus, Arens conducted no statistics at all to test his assumptions, and Altmann (1983) tested the goodness of the non-linear Menzerathian model with F-tests which are very likely to result in misleading interpretations in case of large sample sizes, typical for linguistic data. In fact, as a re-analysis of Arens' data shows, fitting equation (1b) results in a rather poor fit ( $R^2 = 0.70$ ), which is far from being convincing, and consequently sheds doubt on the adequacy of the Menzerathian interpretation.

In an attempt to find some explanation for this poor result by way of a systematic re-analysis of the sentence length – word length problem, Grzybek and Stadlober (2007) and Grzybek et al. (2007) have pointed out a number of possible problems coming into play:

1. *Data Sparsity.* Both the Menzerath-Altmann Law and Arens Law as a special case of it are what one might term “laws of averages”, consequently demanding for a sufficient amount of data points for averages to be reliable. However, due to the large variance of *SL*, an insufficient amount of observations may be available for quite a number of data points of the independent variable. As a consequence, the frequency of observations for each data point has to be guaranteed to prevent random results. In fact, by pooling data into specific classes (as is usual in *SL* analyses), Grzybek, Kelih & Stadlober (2007) arrived at values of  $0.93 \leq R^2 \leq 0.97$ , differences depending on the pooling procedure chosen.
2. *Data homogeneity and text typology.* Given the fact that Arens' original data were based on German literary texts only, the question arises in how far the conclusions made can be generalized and transferred to other text types, as well. Thus, enlarging Arens' text data base by adding literary and scientific prose texts, previously analyzed by Fucks (1955), Grzybek and Stadlober (2007) found the  $R^2$  value to become significantly worse.
3. *Intra-textual vs. Inter-textual approach.* The initial idea of the Menzerath-Altmann Law has been to describe the relation between the constituting components of a given construct and this construct; consequently, the Menzerath-Altmann Law originally was designed in terms of an intra-textual law, relevant for the internal structure of a given text sample. Arens' data, however, are of a different kind, implying inter-textual relations, based on the calculation of sentence length and word length means ( $m_{SL}$ ,  $m_{WL}$ ) for each individual text sample, thus resulting in a vector of arithmetic means. Therefore, in their systematic analysis of 199 Russian texts, Grzybek et al. (2007) obeyed the need to clearly keep the intra-textual and inter-textual perspectives apart. Concentrating on the inter-textual level only, they conducted separate analyses for six different text types, on the one hand, and corpus analyses for the combined data. As a result, they found only very weak evidence on support of Arens Law on an inter-textual level: for the individual text types, the results were between  $0.02 \leq R^2 \leq 0.26$ , for the complete corpus they obtained  $R^2 = 0.49$ . This result coincides with previous observations that obviously, average word length is relatively stable within a given text type – and it is a matter of fact that there can be no variation of word length depending on varying sentence length, if the dependent variable word length displays only poor variation.

## 3. The intra-textual perspective

The present study concentrates on an analysis of the sentence length – word length relation from an intra-textual perspective. Table 1 represents the text data with relevant characteristics.

Table 1  
Text corpus and sub-corpora

Text type	Author	Number of texts	Words		Sentences	
			abs.	rel.	abs.	rel.
Drama	A.P. Čechov	44	67 430	0.28	11125	0.47
Private letters	(various)	120	56 751	0.23	4178	0.18
Literary prose	L.N. Tolstoj	69	74 708	0.31	5 680	0.24
Comments	(various)	60	43 263	0.18	2 556	0.11
Corpus		293	242 152	1.00	23 539	1.00

As can easily be seen, the proportions of sentences and words clearly differ for the different text types; consequently,  $m_{SL}$  and  $m_{WL}$  significantly differ across text types, as has well been documented elsewhere. With this in mind, it will be interesting to analyze the sentence length – word length relation separately for each text type; yet, by way of a first approximation, Fig. 1 offers an overview for the whole corpus.

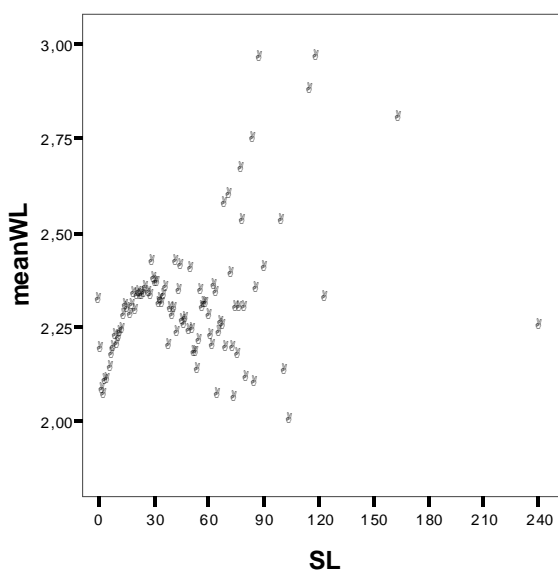


Fig. 1. Word length vs. sentence length: Total Corpus

An inspection of Figure 1 immediately shows the extreme variance of  $m_{WL}$  for long sentences with  $SL \gtrsim 30$ . It is well possible that we are concerned here with linguistic reasons, possibly coming into play; this possibility will be discussed in more detail below. Yet, another possibility must be checked first, which is of statistical rather than linguistic nature. In principle, this reason would concern short sentences as well as long sentences, but particularly long sentences, with  $SL \gtrsim 30$ , are likely to occur relatively rarely. So for a given  $SL$ ,  $m_{WL}$  may be based on a few observations, only, causing a greater variation of  $m_{WL}$ . The increase of word length variation for sentences (and the resulting “loss” of a possibly existing systematic tendency in the  $WL$ - $SL$  relation) might therefore be motivated by merely statistical reasons.

Figures 2 display the frequencies of particular  $SL$  occurrences; indeed, it can easily be

seen, that for all four text types, it is just around  $SL \approx 30$  that the frequency of sentences with the given length decreases to less than 30 observations per class.

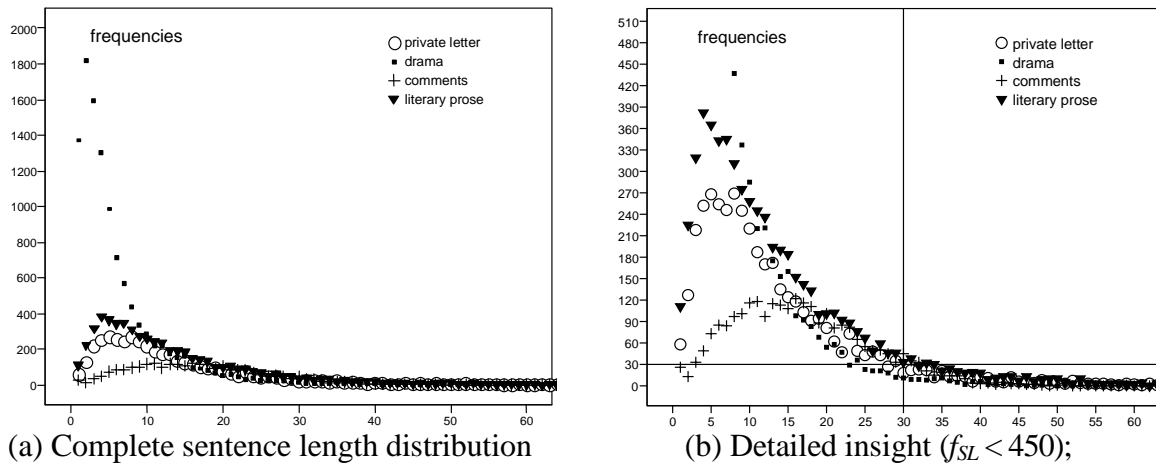


Fig. 2. Sentence length distributions for four text-types

As a consequence, we exclude all occurrences with rare data points for  $m_{WL}$ , by way of an empirical rule of thumb, thus including only data where  $m_{WL}$  is based on 30 observations or more ( $f_{SL} \geq 30$ ); we apply no pooling procedures for the remaining data with less observations, since the type of pooling may be an additional factor influencing the overall result.

Under these circumstances, guaranteeing the postulated minimum of 30 occurrences, a closer look at Figure 3 allows for a more detailed analysis of the overall trend of the core data structure.

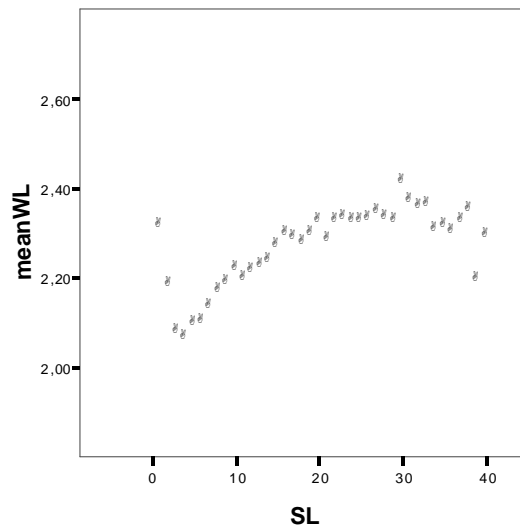


Fig. 3. Word length vs. sentence length: Restricted conditions

Generally speaking, one can now indeed observe a major tendency for longer sentences to be composed of longer words, as predicted by the hypothesis. Yet, there are two important deviations from this overall trend, characterized by two critical points:

1. In very short sentences, the  $SL-WL$  length relation seems to be differently organized as

compared to the bulk of data points: short sentences show a clear decline to a local minimum (in case of the complete corpus, at  $SL = 4$ ), which shall be termed *lower critical point (LCP)*, here. It goes without saying that, for other data material (particularly from other languages), this initial decreasing trend need not be obligatory, and the *LCP* may well be  $LCP \neq 4$ . Anyway, it seems reasonable to assume that we are concerned here with linguistic reasons for this tendency: obviously, very short sentences have no hyposyntactic sub-division and, as a consequence, do not ask for any inter-level Menzerathian control. A detailed analysis of these short sentences must be left for a separate analysis, particularly including *WL* frequency distributions for each of the *SL* classes. In future, it would be desirable to have a common model for all (short and long) sentences; yet, by way of a first approach, we exclude these short sentences from the present study, in order to better concentrate on the bulk of the material, hoping to grasp the general tendency by this procedure.

2. Whereas for sentences with  $4 < SL < 30$ , there seems indeed to be a general tendency for longer sentences to be composed of longer words (as predicted by the hypothesis), there seems to be an *upper critical point (UCP)* for longer sentences with  $SL \gtrsim 30$ . This point is clearly marked by the definite increase of word length variation for these sentences (cf. Figure 3), even after exclusion of occurrences with  $f_{SL} < 30$ . A detailed analysis of this phenomenon goes beyond the scope of this paper; yet, two alternatives lend themselves to interpretation:
  - a. it is possible, that a minimum of  $f_{SL} = 30$  is not sufficient for an average to become stable enough; in this case, we are still concerned with a statistical interpretation of the observed phenomenon,
  - b. it does not seem unlikely that we are concerned here with a (psycho)-linguistically, rather than statistically motivated upper critical point (*UCP*): taking into account human processing limits, Miller's magical rule of  $7 \pm 2$  (and Yngve's linguistic interpretation of it) might well hold true for clause length, and serve as a limitation of the length of clauses or phrases, and, as a consequence, of sentences. Thus, given an average clause length of 5-6 words per clause, the upper limit of information processing on this level might be reached, as a result "de-activating" the Menzerathian control.

In any case, in order to concentrate on the bulk of the material, thus hoping to obtain reliable information on the core of the data structure and grasp its overall tendency, we introduce three empirically motivated restrictions in this study :

- (a)  $f_{SL} > 30$ ,
- (b)  $m_{WL} > LCP$ , and
- (c)  $SL < 30$ .

With these empirical restrictions, it will now be interesting to look not only at the total corpus, but also at the specifics of each of the four different text types. Some basic characteristics of the relevant core data structures are represented in Table 2:

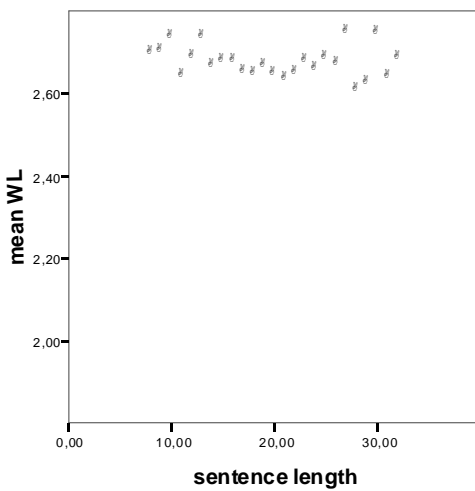
1. The Lower Critical Point (*LCP*) is defined as the minimal  $m_{WL}$  point subsequent to which there is a monotonous increase;
2. the Upper Critical Point (*UCP*) is determined by the empirical restriction of  $f_{SL} > 30$ ;
3. the proportion (in %) of sentences is the percentage of data material representing the core data structure in the interval [*LCP*, *UCP*].

Table 2  
Text corpus and sub-corpora

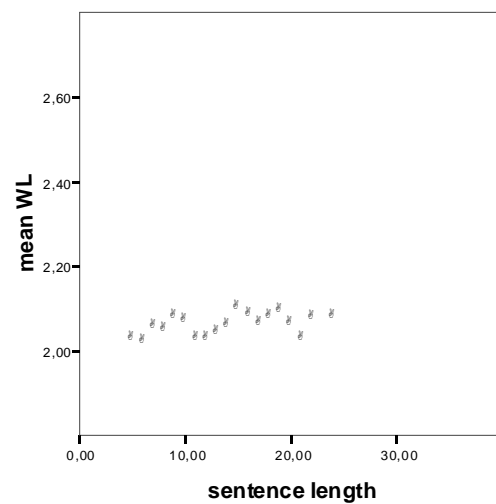
Text type	<i>LCP</i>	<i>UCP</i>	%
Drama	4	22	95.64
Private letters	3	27	90.45
Comments	7	32	94.20
Literary prose	2	31	93.30
<hr style="border-top: 1px dashed black;"/>			
Total	4	40	97.90

As can be seen, both *LCP* and *UCP* differ for the individual text types: Whereas the *LCP* ranges from  $2 \leq LCP \leq 4$ , the *UCP* ranges from  $22 \leq UCP \leq 32$  (in case of the total corpus even reaching  $UCP = 40$ ).

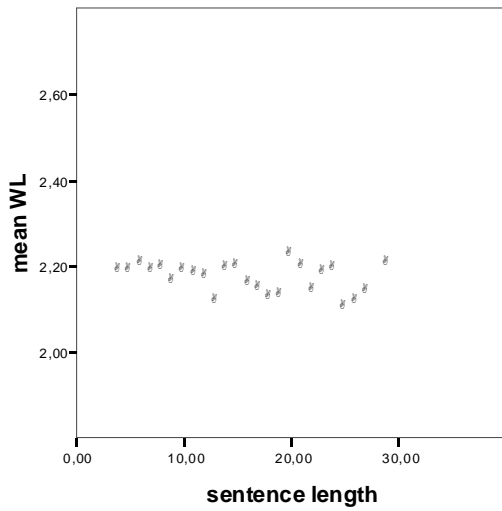
The core data structures for the four text types are represented in Figures 4a-d. With regard to the *SL*–*WL* relation, the results are extremely surprising: quite opposite to expectation, there is almost no increase in  $m_{WL}$  for three of the four text types: rather, in case of the comments, private letters, and dramatic texts,  $m_{WL}$  is almost stable across different *SL* classes. Only for the literary texts, we obtain a convincing fit of  $R^2 = 0.88$  for the non-linear Menzerathian model, with parameter values  $a = 1.93$  and  $b = 0.05$ .



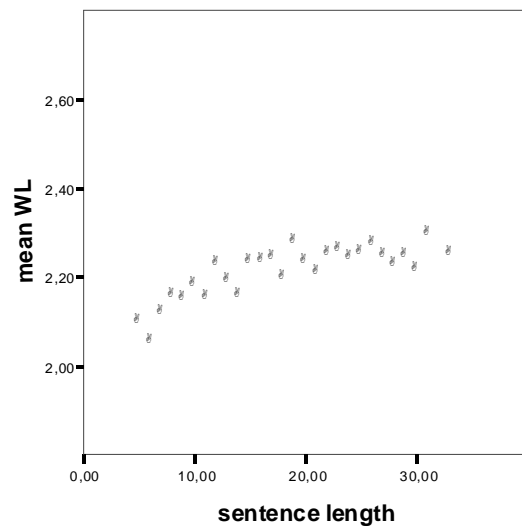
(a) Comments



(b) Drama



(c) Private letters

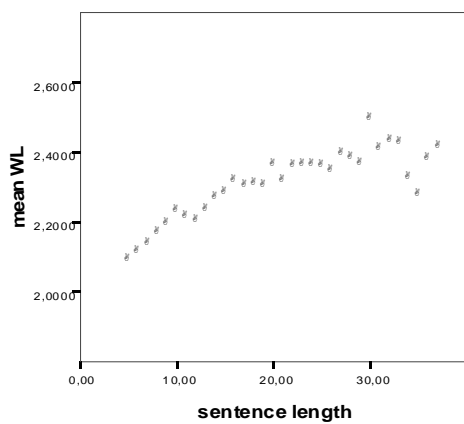


(d) Literary prose

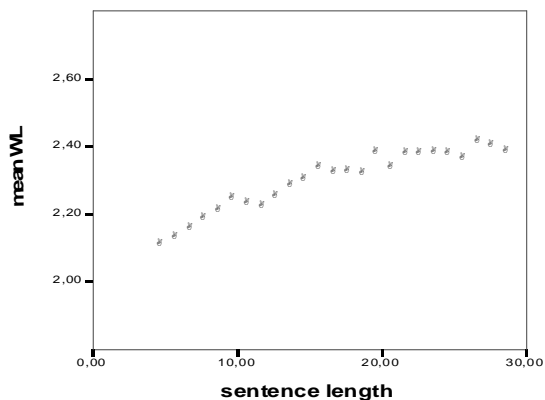
**Fig. 4.** Word length vs. sentence length

In an attempt to find an interpretation of these findings, it seems reasonable to exclude any possible influence of the literary prose texts on the overall corpus. The easiest way to do this, is an additional analysis of a corpus consisting of all comments, private letters, and drama texts, but without the literary texts. This corpus of 167,444 words and 17,859 sentences contains 69.15% of the words and 75.87% of the sentences of the total corpus; its critical points are  $LCP = 4$  (with  $m_{WL} = 2.07$  at this point), and  $UCP = 37$  ( $m_{WL} = 2.42$ ).

Figure 5 (a) shows the  $SL-WL$  tendency for this particular corpus; again, like in the total corpus, there is a fluctuation of  $m_{WL}$  for  $SL > 30$ . Again discarding all sentences with  $SL > 30$ , however, the corpus of comments, drama texts and private letters, with  $R^2 = 0.87$  ( $a = 1.88$ ,  $b = 0.07$ ), shows an almost identical tendency as the literary texts.



(a) Total corpus without literary texts



(b) Total corpus: Core data structure

**Fig. 5.** Word length vs. sentence length



Given these results for the partial corpus (without literary texts), let us now compare them to those for the total corpus. Again, concentrating on the core data structure of the total corpus, excluding short sentences, and cutting off the data at  $SL = 30$ , yields a convincing fit of the Menzerathian non-linear curve: with a determination coefficient of  $R^2 = 0.96$ . Interestingly enough, the parameter values  $a = 1.88$  and  $b = 0.07$  are almost identical with the one obtained for the corpus without the literary prose texts. Figure 5(b) illustrates the overall result.

We thus obtain a number of interesting results:

1. For three of the four analyzed text types (drama, comment, letters), no Menzerathian tendency can be confirmed; only for literary texts, a Menzerathian tendency (Arens Law) can be confirmed;
2. For a partial corpus consisting of these three text types, a Menzerathian (Arens Law) tendency can be confirmed; the same holds true for the total corpus of all four text types.

In attempting to find an answer to the alleged contradictions, it seems reasonable to pay attention to the obviously important factor of data heterogeneity: in case of the partial and total corpora, we are concerned with different text types, each characterized by specific WL and SL characteristics: thus, for the drama texts, we have  $m_{WL} = 2.04$  and  $m_{SL} = 6.06$ , for the letters  $m_{WL} = 2.19$  and  $m_{SL} = 3.58$ , and for the comments  $m_{WL} = 2.67$  and  $m_{SL} = 16.93$ . Only taken together, merged into one common corpus of heterogeneous data, the Menzerathian tendency (Arens Law) appears to be at work. Let us term this phenomenon, which must be subjected to more empirical testing in future, *external textual heterogeneity*.

If this interpretation holds true, a similar hypothesis might be brought forth with regard to the literary texts, as well: in this case, it might well be possible that we are concerned with some kind of *internal textual heterogeneity*, literary texts characteristically being composed of dialogues, descriptive passages, narrative sequences, etc., all of which may well be shaped by different WL and SL characteristics.

Seen from this point, the emergence of the Menzerathian tendency (Arens Law) would have to be interpreted in terms of an index heterogeneity, at least as far as the external perspective is concerned – as to the internal perspective, only some rudimentary insights could be gained in this paper, and more systematic study is necessary in future.

#### 4. Conclusion

The present study offers some important conclusions as to an interpretation of the SL-WL relation along the Altmann-Menzerathian line. Obviously, it seems to work, in case a number of pre-conditions are fulfilled:

- *Minimal sentence length*. For very short sentences ( $SL < 4$ ), the Menzerathian tendency does not seem to play a crucial role; it seems reasonable that this circumstance is motivated by linguistic reasons only, sentences of this length not being subdivided into linguistic sub-units; it goes without saying that the resulting LCP may well be different (or even non-existing) for other languages.
- *Maximal sentence length*. For very long sentences ( $SL > 30$ ), the Menzerathian tendency does not seem to play a crucial role; (psycho)linguistic reasons might be responsible for this circumstance, sentence regulation being at work only as long as a sub-division into sub-units of sentences can be cognitively controlled.
- *Minimal frequency*. Here, we are concerned with a predominantly statistical constraint: if there are not enough (SL) data points as a basis of  $m_{WL}$ , variance is too large to result in some kind of general tendency; accidentally, the UCP of SL around 30 coincides in

most of the data analyzed in this paper with the one explained by maximal SL.

- *Textual heterogeneity*. The Menzerathian principle seems to be of relevance for the *SL-WL* relation only in case sufficient linguistic heterogeneity is guaranteed: as long as the data material to be analyzed consists of homogenous texts (i.e., from a specific text type), *WL* seems to be regulated and, in fact, dominated, by this text type's specific *WL* organization. Only in case data from different text types are combined, the necessary textual heterogeneity is provided for the Menzerathian principle to come into play. It may well be that a literary text as a whole is characterized by this intrinsic heterogeneity, being composed of (homogeneous) text elements such as dialogues, descriptive and narrative sequences, auctorial comments, etc. This might be an explanation why the Menzerathian tendency can be observed in literary texts. It would be particularly interesting to see whether within literary texts, such homogeneous text elements can be isolated which, taken in isolation, do not display any Menzerathian tendencies, yet would, combined into a (heterogeneous) whole. A systematic test of this hypothesis must be left for future research, however.

In addition to these detailed problems, another open question is, if and how very short sentences on the one hand, and long sentences, on the other, can be integrated into one complex model. In other words: It will be an important future task to study (a) in how far the extreme ranges of word and sentence length are characterized by a diverging tendency as compared to the core data structure, and (b) if, both possibly heterogeneous tendencies can yet be incorporated into one overall model. Furthermore, the question of intrinsic heterogeneity, obviously characterizing literary texts, must be subjected to detailed analyses.

## References

- Altmann, G.** (1980). Prolegomena to Menzerath's Law. *Glottometrika 2*, 1–10. Bochum: Brockmeyer,
- Altmann, G.** (1983). H. Arens' «Verborgene Ordnung» und das Menzerathsche Gesetz. In: M. Faust et al. (Eds.), *Allgemeine Sprachwissenschaft, Sprachtypologie und Textlinguistik*: 31-39. Tübingen: Narr.
- Altmann, G., Schwibbe, M.H.** (1989). *Das Menzerathsche Gesetz in informationsverarbeitenden Systemen*. Hildesheim: Olms.
- Antić, G., Stadlober, E., Grzybek, P., and Kelih, E.** (2006). Word length and frequency distributions. In: M. Spiliopoulou et al. (Eds.), *From data and information analysis to knowledge engineering*: 310-318. Berlin: Springer.
- Arens, H.** (1965). *Verborgene Ordnung. Die Beziehungen zwischen Satzlänge und Wortlänge in deutscher Erzählprosa vom Barock bis heute*. Düsseldorf: Pädagogischer Verlag Schwann.
- Fucks, W.** (1955): Unterschied des Prosastils von Dichtern und Schriftstellern. Ein Beispiel mathematischer Stilanalyse. *Sprachforum 1*, 234–241.
- Grzybek, P., Kelih, E., Stadlober, E.** (2007). Long sentences, long words – short sentences, long words? *Presentation at the 31. Jahrestagung der Gesellschaft für Klassifikation: «Data Analysis, Machine Learning, and Application»*. (Freiburg, Germany, March 2007)
- Grzybek, P., Stadlober, E.** (2007). Do we have problems with Arens' law? A new look at the sentence-word relation. In: P. Grzybek and R. Köhler (Eds.), *Exact Methods in the Study of Language and Text*: 205-218. Berlin: de Gruyter.
- Grzybek, P., Stadlober, E., Kelih, E., and Antić, G.** (2005). Quantitative text typology: the impact of word length. In: C. Weihs, and W. Gaul (Eds.), *Classification – The Ubiquitous Challenge*: 53-64. Berlin: Springer.

- Grzybek, P., Stadlober, E., Kelih, E.** (2007). The relationship of word length and sentence length: the inter-textual perspective. In: R. Decker and H.-J. Lenz (Eds.): *Advances in Data Analysis: 611-618*. Berlin: Springer.
- Kelih, E., Grzybek, P., Antić, G., and Stadlober, E.** (2006). Quantitative text typology: the impact of sentence length. In: M. Spiliopoulou et al. (Eds.): *From Data and Information Analysis to Knowledge Engineering: 382-389*. Berlin: Springer.
- Wimmer, G.; Altmann, G.** (2005). Unified derivation of some linguistic laws. In: Köhler, R., Altmann, G., Piotrowski, R. (eds.), *Quantitative Linguistik – Quantitative Linguistics. Ein Internationales Handbuch – An International Handbook: 791-807*. Berlin/New York: de Gruyter.
- Wimmer, G.; Altmann, G.** (2005). Towards a unified derivation of some linguistic laws. In: Grzybek, P. (ed.), *Contributions to the science of text and language. Word length studies and related issues: 329-337*. Dordrecht, NL: Springer.

# **Glottometrics 16**

**Dedicated to Viktor V. Levickij  
on the Occasion of his 70<sup>th</sup> Birthday**

**2008  
RAM-Verlag**

# Glottometrics

**Glottometrics** ist eine unregelmäßig erscheinende Zeitschrift für die quantitative Erforschung von Sprache und Text

**Beiträge** in Deutsch oder Englisch sollten an einen der Herausgeber in einem gängigen Textverarbeitungssystem (vorrangig WORD) geschickt werden

Glottometrics kann aus dem **Internet** heruntergeladen, auf **CD-ROM** (in PDF Format) oder in **Buchform** bestellt werden

**Glottometrics** is a scientific journal for the quantitative research on language and text published at irregular intervals

**Contributions** in English or German written with a common text processing system (preferably WORD) should be sent to one of the editors

Glottometrics can be downloaded from the **Internet**, obtained on **CD-ROM** (in PDF) or in form of **printed copies**

## Herausgeber – Editors

G. Altmann	ram-verlag@t-online.de
K.-H. Best	kbest@gwdg.de
F. Fan	fanfengxiang@yahoo.com
P. Grzybek	peter.grzybek@uni-graz.at
L. Hřebíček	ludek. <a href="mailto:hrebicek@seznam.cz">hrebicek@seznam.cz</a>
R. Köhler	koehler@uni-trier.de
J. Mačutek	jmacutek@yahoo.com
G. Wimmer	wimmer@mat.savba.sk
A. Ziegler	arne.ziegler@uni-graz.at

**Bestellungen** der CD-ROM oder der gedruckten Form sind zu richten an  
**Orders** for CD-ROM's or printed copies to

RAM-Verlag      [RAM-Verlag@t-online.de](mailto:RAM-Verlag@t-online.de)

**Herunterladen / Downloading:**      <http://www.ram-verlag.de>

Die Deutsche Bibliothek – CIP-Einheitsaufnahme

Glottometrics. –16 (2008-03-136) –. – Lüdenscheid: RAM-Verl., 2008  
Erscheint unregelmäßig. – Auch im Internet als elektronische Ressource  
unter der Adresse <http://www.ram-verlag.de> verfügbar.-

Bibliographische Deskription nach 16 (2008)

**ISSN 1617-8351**

# Contents

<b>Ján Mačutek</b> Runes: complexity and distinctivity	1-16
<b>Nadia Yesypenko</b> Writer's voice in the texts of “Peace and War” themes	17-26
<b>Karl-Heinz Best</b> Word length in Persian	27-30
<b>Ioan-Iovitz Popescu, Gabriel Altmann</b> Zipf’s mean and language typology	31-37
<b>Nora Heinicke</b> Wortlängenverteilungen in französischen Briefen eines Autors	38-45
<b>Emmerich Kelih</b> Modelling polysemy in different languages: A continuous approach	46-56
<b>Marina Knaus</b> Zur Verteilung rhythmischer Einheiten in russischer Prosa	57-62
<b>Solomija Buk, Ján Mačutek, Andrij Rovenchak</b> Some properties of the Ukrainian writing system	63-79
<b>Zuzana Martináková, Ioan-Iovitz Popescu, Ján Mačutek, Gabriel Altmann</b> Some problems of musical texts	80-110
<b>Peter Grzybek, Emmerich Kelih, Ernst Stadlober</b> The relation between word length and sentence length: an intra-systemic perspective in the core data structure	111-121
<b>History of Quantitative Linguistics</b>	122-131
<b>Karl-Heinz Best</b> XXXII. Helmut Meier (1897-1973)	122-124
<b>Karl-Heinz Best</b> XXXIII. Adolf Busemann (1887-1967)	124-127
<b>Karl-Heinz Best</b> XXXIV. Kaj Brynolf Lindgren (1822-2007)	127-131