

Text difficulty and the Arens-Altman law

Peter Grzybek

1 Introduction

The study of text difficulty is considered to be an important issue for many branches of applied research. In the fields of journalism or education, for example, it is particularly important to know if (or to what degree) a given text is likely to cause difficulties for a recipient, or a group of recipients, i.e., if it is likely to be on an adequate (intended) level of difficulty or beyond.

In order to achieve this goal, a specific line of text difficulty research has developed over the last decades, beginning in the 1920s, which attempts to combine linguistic analysis with informants' ratings of text difficulty.¹ Text difficulty thus is a double-faced kind of empirical research in two directions, either of which may be emphasized in individual studies: it is text-based, on the one hand, and informant-oriented, on the other. Due to this dual perspective, alternative terms such as 'text readability'² or 'text comprehensibility' have been used to refer to the related area(s) of research, the first term emphasizing the predominantly written (rather than oral) basis of communication, the second being broader in its understanding. Compared to these alternatives, 'text difficulty' as a term primarily refers to the analysis of linguistic structures, aiming at the identification and characterization of linguistic factors rendering a given text more or less easily comprehensible to a given person (or a group of persons), and at the (cor)relation of these structures to informants' ratings about text difficulty. Such a definition is in line with research from the last decades: Klare (1963: 1), for example, understands this term as referring to "the ease of understanding or comprehension due to the style of writing", and DuBay (2004: 3), more recently, has defined the overall aim of text difficulty research as the study of "what makes some texts easier to read than others".

Given this general orientation, research in this field, from its beginnings on, has continually tried to develop, modify and improve formulae to predict text difficulty and, by way of it, prognose comprehension ability. This is to say that attempts have been undertaken to develop measures of text difficulty, including formulae which combine quantitative (or quantified) linguistic characteristics in such a way that these characteristics serve as (possibly combined and

1. Informants may be either be recipients, mainly readers, or experts in the given field, such as teachers, librarians, publishers, lecturers of publishing houses, etc.

2. 'Text readability' in turn should not be confused with 'text legibility' which concerns factors such as typeface and layout of texts.

specifically weighted) factors for an optimized prediction of text difficulty. In the history of research³ starting in the early 1920s, a number of relevant phases can be distinguished, in which researchers have tried to identify linguistic factors to be good indicators and predictors of text difficulty⁴. Early work as e.g. by Lively and Pressey (1923) mainly concentrated on lexical analysis; here, two major approaches can be distinguished: research concentrated on either the (relative) number of different words in a given text⁵, or on references to frequency lists.⁶ Subsequent work attempted to enlarge the linguistic spectrum and identify further factors, guided by the principle «The more, the better»: thus, authors like Gray and Leary (1935) already used a collection of 64 linguistic variables. Later, possible interactions between different linguistic factors became focused, in order to arrive at higher levels of correlation between attributed text difficulty and the combination of a set of linguistic variables. In this direction, two important results were obtained: first, many linguistic variables were highly intercorrelated, and second, an increase of the number of linguistic variables did not generally raise the correlation coefficient. Since, therefore, the use of more variables may be only minutely more accurate, but much more difficult to measure and apply, the next step included the reduction of variables and the identification of maximally predictive factors.

As a consequence, many different formulae were developed over the following years; Klare (1981) noted there were over 200 published formulae to measure text difficulty. All of these formulae have been developed by inductive-empirical approaches, typical for research in this field. Most of these formulae differ less as to the linguistic factors included, rather than how they are weighted. Among those factors re-occurring most frequently in all these formulae, are factors such as word frequency, amount of different words, average sentence length, average word length, and others (cf. Amstad 1978: 48f.).

From the perspective of quantitative linguistics in general, and synergetic linguistics, in detail, the high degree of relatedness between the various linguistic factors is not surprising; after all, it is well-known that both frequency and length characteristics of linguistic units on all analytical levels are closely

-
3. Since there are a number of informative surveys on this topic, this need not be presented here in detail.
 4. Klare (1963: 4), for example, has distinguished between four phases of development: according to him, the early 'pioneer phase' (1921–1934) was followed by the development of detailed (1934–1938), efficient (1938–1953) and specialized (1953ff.) formulae.
 5. This approach is well-known today as the study of 'lexical richness', usually including some kind of lexical type-token ratio. As we know today, there are quite a number of theoretical problems with this approach as, e.g., the dependence of the type-token ratio on text length. Additionally, it should be mentioned that in these early studies, no specific definition of 'word' has been used and, as a consequence, no distinction between 'word' and 'word form' (or lemma) has been made.
 6. The early studies were mainly based on E.L. Thorndike's (1921, 1932), or Thorndike's and Lorge's (1944) lexical frequency analyses; later studies rather referred to G.K. Zipf's works as a reference line, which are better known today.

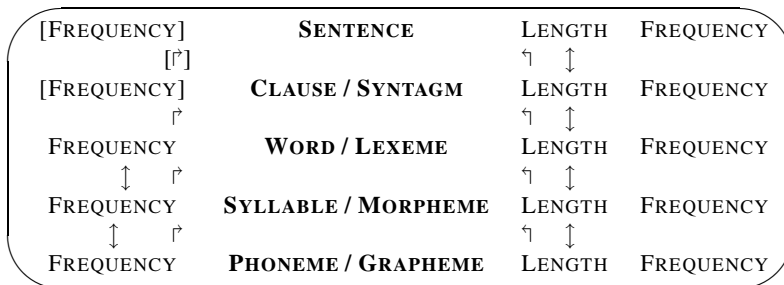
related and mutually interwoven. As a consequence, it is almost self-evident that if text difficulty is to be measured by reference to linguistic characteristics, it is sufficient to concentrate on only a few factors.

In this respect, the Flesch Reading Ease Index (*REI*), developed by Flesch (1948) with regard to English texts, is probably the most quoted and one of the easiest to apply. It is the result of a “simple” linear regression, i.e. combination of the average word length (*WoL*) and average sentence length (*SeL*) of a given text as the only two relevant factors (in addition to a constant):

$$REI_{engl} = 206.835 - (1.015 \cdot SeL) - (84.6 \cdot WoL) \quad (1)$$

Although quite simple at first sight, this formula is still today considered to be very efficient⁷ and probably it is just due to its easy application that it is continuing to be one of the most widely used to measure text difficulty. Last not least, it is just the ambivalence between simplicity and efficiency of this formula which has given rise to skepticism, partly motivated by the lack of the formula’s theoretical foundation. In this context, the validity of this formula has been generally called into question emphasizing the fact that “isolated linguistic units” are no adequate means for measuring text difficulty.

This view contradicts, of course, the above-mentioned synergetic interrelations between linguistic units, the relevance of which for text difficulty research have hardly ever been theoretically reflected in the whole research area. Therefore Best (2006), in his critical analysis of this discussion, is fully correct in objecting and countering that there are no isolated units in language. Particularly the word may be seen in the center of ‘horizontal’ and ‘vertical’ interrelations; as is well-documented, the word is part of a complex control circuit, the most basic factors of which are word length, semantic complexity, co-textuality, and word frequency (cf. Köhler and Altmann 1986: 261). Other relevant elements of this self-regulating dynamic system are syllable/morpheme length, clause length, sentence length, etc., and their respective frequencies. The following schema illustrates some basic synergetic processes; it makes clear that frequency and length characteristics of linguistic units stand in close self-regulating relations:



7. In comparative studies, the Flesch formula has repeatedly turned out to be the most efficient of those which need no word list (cf. Amstad 1978: 64).

As a result, Best (2006) correctly concludes: “Readability formulae, based on sentence and word length, indirectly measure substantially more than is expressed in these formulae, due to the manifold interactions between linguistic units.” This view contains, of course, no theoretical foundation as to the question which specific factors influence text difficulty in what way or to what degree; yet it offers a theoretically based post-hoc answer to the question why the reduction to only a couple of seemingly elementary factors has made this concept to have such a success story.

Notwithstanding this insight, there is a whole bunch of crucial questions which continue to be unsolved. A major problem is the language-specific character of Flesch’s *REI*: as was pointed out above, formula (1) was originally developed for English texts in the late 1940s. In later attempts to apply this formula to other languages, it soon turned out that language-specific adaptations were necessary, mainly due to the interest of having results on a scale from 0 to 100 in each language. Thus, for example, for Dutch, French, Spanish, German and Ukrainian the following adaptations were suggested⁸, all following the general expression $REI = C - a \cdot WoL - b \cdot SeL$:

$$REI_{dutch} = 195 - (0.66 \cdot WoL) - (2 \cdot SeL), \quad (1a)$$

$$REI_{french} = 207 - (73.6 \cdot WoL) - (1.015 \cdot SeL), \quad (1b)$$

$$REI_{german} = 180 - (58.5 \cdot WoL) - SeL, \quad (1c)$$

$$REI_{spanish} = 206.84 - (77 \cdot WoL) - (0.93 \cdot SeL), \quad (1d)$$

$$REI_{ukrainian} = 206.84 - (28.3 \cdot WoL) - (5.93 \cdot SeL). \quad (1e)$$

As can be seen, the language-specific differences between these formulae consist in different weights for *WoL* and *SeL*, i.e. in different parameter values for *a* and *b*. *WoL* and *SeL* thus represent two crucial factors in measuring text difficulty across languages; yet, either their importance as separate factors, or their specific interrelation (i.e., the relation between *WoL* and *SeL*), clearly differs for individual languages.

Unfortunately, no systematic cross-linguistic studies are available which might explain what causes, or motivates, the observed differences in weighting. From a theoretical perspective, Best’s (2006) reference to the synergetic specifics of language offers a good starting point for research in this direction. In this context, particularly the *WoL* – *SeL* relation has recently been studied in detail, both from an inter-textual and intra-textual perspective; whereas the first concentrates on relations within a given text (or groups of texts), the second compares more than one textual object and studies the relation between them. For both perspectives, law-like regularities have been postulated and demon-

8. Cf. Kandel and Moles (1958), Fernández Huerta (1959), Brouwer (1963), Amstad (1978), Partiko (2001: 257).

strated to exist. From an *intra-textual perspective*, we are concerned with the *Menzerath-Altman law*, relevant for the relation between a given construct and its constituting components within a given text (notwithstanding the possibly intervening level of clauses coming into play, on an intermediary level between sentence and word). As compared to this, the *inter-textual relation* is covered by the *Arens-Altman law*, based on the calculation of the mean length of words (\bar{x}) and sentences (\bar{y}) in a series of text samples, resulting in two vectors of arithmetic means ($\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$).

In order to gain insight into the specific role *SeL* and *WoL* play for text difficulty in the individual languages, it seems reasonable, therefore, to study relevant data on the background of the Arens-Altman law. Since such a systematical approach has never been undertaken before, a first approach into this direction should start with one language only. But even with this restricting focus, it is of utmost importance to pay due attention to yet another circumstance: as recent analyses have shown (Grzybek et al. 2007, Grzybek and Stadlober 2007, Grzybek et al. 2008), both *WoL* and *SeL* are not constant within a given language (i.e., are not ‘typical’ of a given language as a whole); rather, they differ for specific discourse types within a language. It seems likely that this finding is also relevant for the *WoL* – *SeL* relation, but this possibility, too, has never been submitted to systematical reflection.

In the following analyses, these objectives shall be pursued, using German language material, strictly controlling text type. Since the perspective should be cross-linguistic right from the beginning, it seems reasonable to immediately provide a meta basis adequate for comparison. In this respect, suggestions developed by Estonian scholar Tuldava in a series of articles (1993a,b), turn out to be of utmost importance, since they contain a language-independent formula of measuring text difficulty (*TD*), also based on *WoL* and *SeL*, only:

$$TD = WoL \cdot \ln(SeL) . \quad (2)$$

This formula has remained rather unknown in the field of text difficulty research. As a consequence, its efficiency has never been generally tested; specifically, no systematic comparisons with Flesch’s *REI* or any one of its language-specific adaptations have ever been undertaken. Tuldava himself applied his formula (2) to a sample of 20 German texts of different types (text books, journalistic, literary prose, scientific). Comparing the results obtained to Flesch’s original *REI* formula (1), rather than to Amstad’s German adaptation (1c), Tuldava (1993a: 78) found a close rank correlation of $\varphi = -0.97$ between these two measures.⁹ Tuldava did not attempt to establish a detailed regression equation, which would allow for the transformation of one measure to the other.

9. As a re-analysis of his data shows, this correlation is highly significant ($p < 0.001$), with the linear regression $TD^* = 7.16 - 0.051 \cdot REI$.

If this result were confirmed on a broader basis of linguistic material, this would mean that Tuldava's formula (2) could indeed serve as a basis for cross-linguistic comparisons, for which no language-specific parameter estimations would be needed. More importantly, this would be an important step in the direction outlined above, both in practical and theoretical respects:

- From a *practical* point of view, the application of Tuldava's parameter-free formula would not only imply the option of measuring text difficulty without knowledge of language-specific parameters (i.e., weights), but, in addition to this, the results obtained might easily be transformed to fit one of the 'established' Flesch measures mentioned above.
- From a *theoretical* perspective, insight might be gained as to the question how *WoL* and *SeL*, either as individual factors or as a complex combination in their self-regulating interrelation, influence text difficulty.

The detailed study of the *WoL* – *SeL* relation is of utmost importance in yet another respect for text difficulty research: If *WoL* can be characterized to depend on *SeL*, as predicted by the Arens-Altman law, then Tuldava's formula (2) might even be further reduced to one linguistic variable, only. At first sight, it might be equally plausible to substitute either the *WoL* or the *SeL* variable by the theoretical value to be expected according to the Arens-Altman law; however, with *WoL* being the dependent variable, rather than *SeL*, it seems more appropriate to substitute the *WoL* variable, the more since the latter displays much less variation than *SeL* in a given text. In fact, the idea to substitute *WoL* has been brought forth by Tuldava (1993a), but it has never been empirically tested, due to insufficient research on the Arens law.

2 Analysis

As to appropriate data serving as material for our study, German texts analysed by Bamberger and Vanecek (1984) in their study on readability of school texts seem to be adequate. The authors investigated the readability of 380 texts from primary and lower secondary level textbooks; in detail, they analyzed 240 special texts [Sachtexte], and 120 literary prose texts for adults (i.e., youth literature). These texts were evaluated by an expert team according to their appropriateness for different school grades, each text being attributed to a particular difficulty level (*DL*). The authors then applied a variety of readability formulae, taking into account a large number of different linguistic factors which were tested for different levels from grades four through twelve. The linguistic characteristics of these factors are not relevant here; for our purposes, it may suffice to say that among others, average values for *WoL* and *SeL* were calculated for all texts, and these data shall serve for the subsequent re-analysis.

Figures 1a and 1b show the relation between *WoL* and *SeL* for the 380 texts: Figure 1a shows the original data points, in Figure 1b the latter are pooled in

groups of ten each, in order to make the overall tendency appear more transparent. As can be seen, there is an obvious trend of *WoL* to increase with increas-

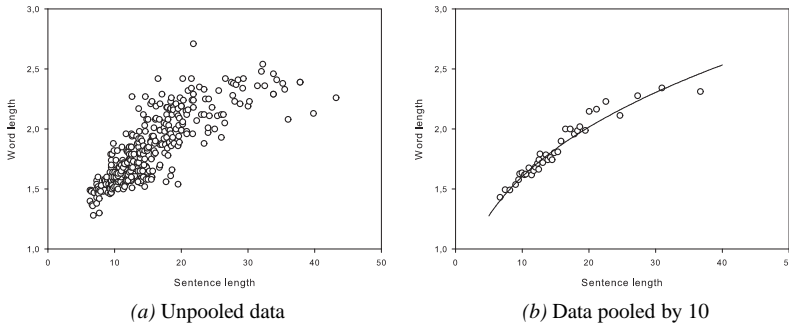


Figure 1: Dependence of *WoL* on *SeL* for 380 German texts from Bamberger and Vanecek (1984)

ing *SeL*; this tendency is particularly clearly expressed in Figure 1b. According to the Arens-Altman law, this relation may be modeled by the function $WoL = a \cdot SeL^b$: in fact, with parameter values $a = 0.75$ and $b = 0.33$, the fit turns out to be very good ($R^2 = 0.95$), as can also be seen from the regression curve added in Figure 1b.

These findings are in accordance with the Arens-Altman law and the hitherto undoubted assumption that, within a given language, the *WoL-SeL* relation on the inter-textual level can be modeled without distinction of text types. However, extending the data base of 380 texts by adding the above-mentioned 117 data sets from the original Arens (1965) study, analogically pooled by items of ten each, radically changes this view. Figure 2 clearly shows that the literary prose texts studied by Arens display the same overall trend of *WoL* increasing with an increase of *SeL*, but in a different way as compared to the schoolbook texts. This finding asks for a differentiated analysis of all three text types separately.

Figure 2 shows the resulting tendencies in detail: Quite obviously, there is an increase of *WoL* with an increase of *SeL* for all three text types.

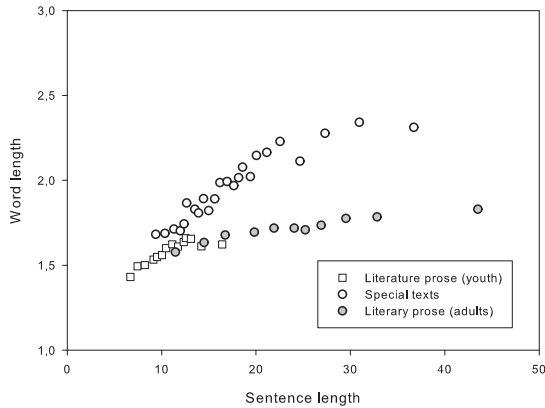


Figure 2: Dependence of *WoL* on *SeL* for 497 German texts (data pooled by 10)

Yet, the kind of increase differs for each of them; this fact is corroborated by the divergent parameter values, which are represented in Table 1.¹⁰

Table 1: Fitting results for three text types

Texts	N	a	b	R^2
Youth literature	140	0.89	0.24	0.9956
Adult literature	117	1.23	0.11	0.9658
Special texts	240	0.69	0.37	0.9562

Summarizing, we can say that no simple substitution of either the *WoL* or the *SeL* variable is possible for Tuldava's *TD* formula, since the relation between *SeL* and *WoL* is not constant within a given language, but differs for text types. It is a task for future research to find out which and how many text types must be distinguished in this respect; it seems to be reasonable, however, to assume that we are concerned with the same kind of discourse types which have been identified to be relevant for the discrimination of discourse type on the basis of 'simple' *WoL* and *SeL* studies (cf. Grzybek et al. 2005; Kelih et al. 2006).

10. Interestingly enough, youth literature and adult literature seem to follow an identical kind of increase, though at different ends of the regression curve: joining the pooled data points for the 257 literary texts in a common type of 'literature' results in a good fit ($R^2 = 0.92$); in this case, we obtain parameter values for $a = 1.25$ and $b = 0.10$, which come very close to those adult literature. Nevertheless, the two literary text groups shall be treated separately in the subsequent analyses.

3 Text difficulty and text types

With these relations established, we can now come back to the question of text difficulty, separately for each of the two text types (i.e., 120 texts from youth literature and 240 special texts). Figures 3a and 3b present the results for Amstad's and Tuldava's formulae (1c) and (2), respectively.

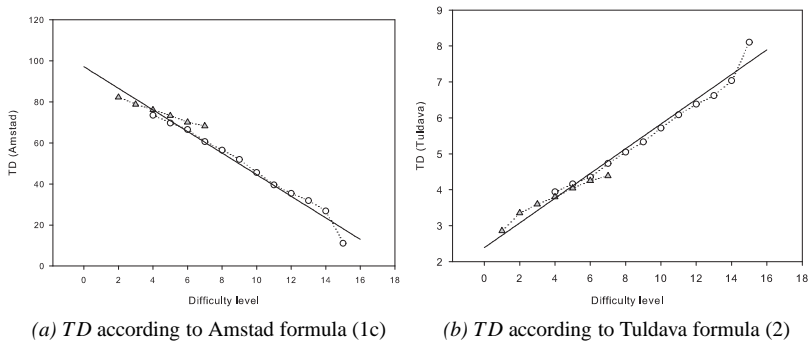


Figure 3: Text difficulty (TD) for 380 German texts (pooled by 20)

An inspection of Figure 3 allows for a number of important observations:

1. As expected, there is a clear major tendency of DL and TD being closely correlated; this tendency holds for both formulae, though with opposite directions. Ignoring text type specifics, the dependency is of clearly linear kind, with a high correlation coefficient of $r = 0.99$ in both cases.
2. Whereas there seem to be clear differences in the kind of relation between word and sentence length for the two text types – at least this was the result of the analyses discussed above (cf. Figure 2) –, the corresponding TD values seem to follow a common tendency (notwithstanding difficulty differences, of course). Obviously, particular text types have their own specific mechanisms of ruling TD , which allows, as a consequence, for a common analytical procedure. As long as no additional data change the picture, or further interpretations are available, it seems reasonable to consider the relation between DL and TD to be linear, across text types as well as within a given text type (with $r > 0.97$ in all four cases). Still, it remains an open question whether or not TD can be reasonably defined without taking into account text typological specifics.
3. Regardless of possible text typological specifics, it turns out that, at least for German, the language-independent measure for TD according

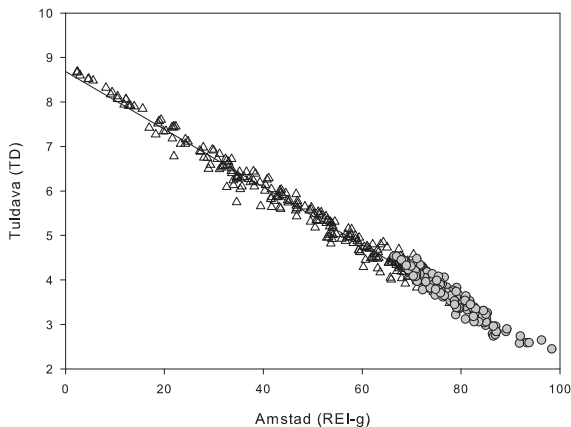


Figure 4: Comparison of Amstad's *REI* and Tuldava's *TD* indices for text difficulty

to Tuldava's formula (2) is equally efficient in predicting *DL* as is Amstad's language-specific adaptation (1c) of Flesch's *REI* to German, the correlation between both measures being highly significant ($r = 0.99$). Figure 4 shows the correlation between both measures, combined for both text types, but with distinct marks. This confirms Tuldava's above-mentioned observations on a broader data basis; additionally, it is based on the specific German adaptation of Flesch's *REI*, rather than on the original developed for English texts. In fact, both formulae turn out to measure in principle the same, though on different scales; as a consequence, they can be transformed one into the other. With regard to the 380 texts analyzed here, for example, the transformation from Tuldava's *TD* value to Amstad's scale might be easily calculated by way of the equation $REI_{german}^* = 133.09 - 15.24 \cdot TD$; alternatively, the transformation from Amstad's scale to Tuldava's value can easily be achieved by calculating $TD^* = 8.68 - 0.065 \cdot REI_{german}$. It goes without saying that, before generally applying these transformations to German texts, more text types must be studied, covering the whole textual spectrum. It is highly probable that this will result in a more or less considerable modification of these transformational procedures; by way of a comparison, the transformation from Flesch's original *REI* into Tuldava's *TD* would result in the equation $TD^* = 6.72 - 0.05 \cdot REI$, which also slightly differs from the figures given in footnote 9.

4 Results and conclusions

A first major result of the present study is the finding that, at least for German, text difficulty can be measured without any language-specific adaptation. As compared to Amstad's German adaptation of Flesch's *REI*, Tuldava's *TD* provides practically the same exactness of predictability, practically without loss of information. This finding is of relevance not only for German: if it can be corroborated for further languages, no language-specific adaptations, and no parameter estimations, will be necessary in future. Tuldava's formula may be considered to be universally valid; but this is a matter of boundary conditions in the individual languages; at present, we have no idea as to this point which represents an interesting linguistic question in its own right, namely, to what extent the formula works in which way (i.e., with which parameters) for which languages.

A second major result is that possibly no text typological specifics need to be taken into account when measuring text difficulty with Tuldava's *TD*: Since word length and sentence length are the only two characteristics taken into consideration in this formula, their interrelation has been submitted to a detailed analysis in this study. This analysis results in the observation that, within a given language, text typological differences do exist, but might not play a crucial role for measuring *TD*; rather, it seems possible that *TD* is the result of a language-intrinsic control mechanism, which allows for the application of a common (unique) procedure in text difficulty analysis.

Given these overall results, a number of important tasks remain to be tackled by future research:

1. As compared to the history of text difficulty research, much more systematic study is necessary; this concerns both cross-linguistic comparisons and intra-lingual specifics of text types:
 - (a) Within a given language, attention must be paid to (the comparability of) different text types; for each of them the specific relation between word and sentence length must be studied.
 - (b) As to cross-linguistic studies, the application of Tuldava's formula and its comparison with language-specific formulae seems to be an extremely promising way; in these inter-lingual comparisons too, of course, due attention must be paid to text typology to compare only like to like.
2. As suggested by Tuldava (1993a), the value for either word length or sentence length may be substituted, theoretically, one for the other. A necessary pre-condition for this substitution is, of course, knowledge about the specific relation between word length and sentence length (be it for a given language, in general, or for specific text groups, in particular). In this respect, it has not been considered sufficiently thus far that, within

a given language, this relation may differ across text types; therefore, before such substitutions, much more systematic study on the *WoL-SeL* relation along the Arens-Altman law and its text type specific boundary conditions is necessary.

3. Tuldava's formula and its efficiency remain almost unexplained; it is obvious that the logarithm included leads to a weight reduction of sentence length, but for the time being, there is no explanation in sight why this weight reduction should be logarithmic. It seems reasonable to assume that controlling the relation between word and sentence length will yield relevant insight into this question, the logarithm possibly turning out to be but a good approximation. In any case, it would be desirable either to strive for a theoretical explanation of the logarithmic weight or to replace the logarithm by a parametric model, the parameters of which, in turn, are then open to be interpreted.

References

- Arens, H.
1965 *Verborgene Ordnung. Die Beziehungen zwischen Satzlänge und Wortlänge in deutscher Erzählprosa vom Barock bis heute.* Düsseldorf: Pädagogischer Verlag Schwann.
- Amstad, T.
1978 *Wie verständlich sind unsere Zeitungen?* Diss., University of Zürich.
- Bamberger, R.; Vanecek, E.
1984 *Lesen – Verstehen – Lernen – Schreiben.* Wien: Jugend und Volk.
- Best, K.-H.
2006 “Sind Wort- und Satzlänge brauchbare Kriterien zur Bestimmung der Lesbarkeit von Texten?” In: Wichter, S.; Busch, A. (eds.), *Wissenstransfer*: Frankfurt/M.: Lang; 21–31.
- Brouwer, R.H.M.
1963 “Onderzoek naar de leesmoeilijkheid van Neerlands proza”, in: *Pedagogische Studiën*, 40; 454–464.
- DuBay, W.H.
2004 *The Principles of Readability.* Costa Mesa, CA: Impact Information.
- Fernández Huerta, J.
1959 “Medidas sencillas de lecturabilidad”, in: *Consigna* 214; 29–32.
- Flesch, R.
1948 “A New Readability Yardstick”, in: *Journal of Applied Psychology*, 32/3; 221–233.
- Gray, W.S.; Leary, B.
1935 “What makes a book readable.” Chicago: Chicago University Press.
- Grzybek, P.; Stadlober, E.
2007 “Do We Have Problems With Arens’ Law? A New Look at the Sentence-Word Relation.” In: Grzybek, P.; Köhler, R. (eds.), *Exact Methods in the Study of Language and Text. Dedicated to Professor Gabriel Altmann on the Occasion of His 75th Birthday.* Berlin / New York: Mouton de Gruyter; 205–218.
- Grzybek, P.; Stadlober, E.; Kelih, E.
2007 “The Relation of Word Length and Sentence Length: The Inter-Textual Perspective.” In: Decker, R.; Lenz, H.-J. (eds.), *Advances in Data Analysis.* Berlin etc.: Springer; 611–618.
- Grzybek, P.; Kelih, E.; Stadlober, E.
2008 “The relation between word length and sentence length. An intra-systemic perspective in the core data structure”, in: *Glottometrics*, 16; 111–121.
- Grzybek, P.; Stadlober, E.; Kelih, E.; Antić, G.
2005 “Quantitative Text Typology: The Impact of Word Length.” In: Weihs, C.; Gaul, W. (eds.), *Classification. The Ubiquitous Challenge.* Heidelberg, New York: Springer, 53–64.

- Kandel, L.; Moles, A.
1958 "Application de l'indice de Flesch à la langue français", in: *Cahiers d'Etudes de Radio-Television*, 19; 253–274.
- Kelih, E.; Grzybek, P.; Antić, G.; Stadlober, E.
2006 "Quantitative Text Typology. The Impact of Sentence Length." In: Spiliopoulou, M.; Kruse, R.; Nürnberger, A.; Borgelt, C.; Gaul, W. (eds.), *From Data and Information Analysis to Knowledge Engineering*. Heidelberg, Berlin: Springer, 382–389.
- Klare, G.R.
1963 *The measurement of readability*. Ames, Iowa: Iowa State University Press.
1981 "Readability indices: do they inform or misinform?", in: *Information design journal* 2; 251–255.
- Köhler, R.; Altmann, G.
1986 "Synergetische Aspekte der Linguistik", in: *Zeitschrift für Sprachwissenschaft*, 5; 253–265.
- Lively, B.A.; Pressey, S.L.
1923 "A method for measuring the 'vocabulary burden' of textbooks", in: *Educational administration and supervision* 9; 389–398.
- Mikk, J.
2000 *Textbook: Research and Writing*. Frankfurt/M. etc.: Lang.
- Partiko, Z.V.
2001 *Zagal'ne redaguvannja. Normativni osnovi*. L'viv: Afiša.
- Thorndike, E.L.
1916 "An improved scale for measuring ability in reading", in: *Teachers college record*, 17; 40–67.
1921 *The teacher's word book*. New York: Bureau of Publications, Teachers College, Columbia University.
1932 *A teacher's word book of 20,000 words*. New York: Bureau of Publications, Teachers College, Columbia University.
- Thorndike, E.L.; Lorge, I.
1944 *The teacher's word book of 30,000 words*. New York: Bureau of Publications, Teachers College, Columbia University.
- Tuldava, J.
1993a "Measuring text difficulty." In: Altmann, G. (ed.), *Glottometrika 14*. Trier: Wissenschaftlicher Verlag wvt; 69–81.
1993b "The statistical structure of a text and its readability." In: Hřebíček, L.; Altmann, G. (eds.), *Quantitative Text Analysis*. Trier: Wissenschaftlicher Verlag wvt; 215–227.

Text and Language

**Structures · Functions · Interrelations
Quantitative Perspectives**

Edited by

**Peter Grzybek
Emmerich Kelih
Ján Mačutek**

Advisory Editor

Eric S. Wheeler

**Praesens Verlag
Wien 2010**

Contents

Preface <i>Peter Grzybek, Emmerich Kelih, Ján Mačutek</i>	vii
Quantitative analysis of Keats' style: genre differences <i>Sergej Andreev</i>	1
Word-length-related parameters of text genres in the Ukrainian language. A pilot study <i>Solomija Buk, Olha Humenchyk, Liliya Mal'tseva, Andrij Rovenchak</i>	13
On the quantitative analysis of verb valency in Czech <i>Radek Āech, Ján Mačutek</i>	21
A link between the number of set phrases in a text and the number of described facts <i>Łukasz Dębowski</i>	31
Modeling word length frequencies by the Singh-Poisson distribution <i>Gordana Đuraš, Ernst Stadlober</i>	37
How do I know if I am right? Checking quantitative hypotheses <i>Sheila Embleton, Dorin Uritescu, Eric S. Wheeler</i>	49
Text difficulty and the Arens-Altman law <i>Peter Grzybek</i>	57
Parameter interpretation of the Menzerath law: evidence from Serbian <i>Emmerich Kelih</i>	71
A syntagmatic approach to automatic text classification. Statistical properties of <i>F</i> - and <i>L</i> -motifs as text characteristics <i>Reinhard Köhler, Sven Naumann</i>	81
Probabilistic reading of Zipf <i>Jan Králík</i>	91
Revisiting Tertullian's authorship of the <i>Passio Perpetuae</i> through quantitative analysis <i>Jerónimo Leal, Giulio Maspero</i>	99
Textual typology and interactions between axes of variation <i>Sylvain Loiseau</i>	109

Rank-frequency distributions: a pitfall to be avoided <i>Ján Mačutek</i>	119
Measuring lexical richness and its harmony <i>Gregory Martynenko</i>	125
Measuring semantic relevance of words in synsets <i>Ivan Obradović, Cvetana Krstev, Duško Vitas</i>	133
Distribution of canonical syllable types in Serbian <i>Ivan Obradović, Aljoša Obuljen, Duško Vitas, Cvetana Krstev, Vanja Radulović</i>	145
Statistical reduction of the feature space of text styles <i>Vasilij V. Poddubnyj, Anastasija S. Kravcova</i>	159
Quantitative properties of the Nko writing system <i>Andrij Rovenchak, Valentin Vydrin</i>	171
Distribution of motifs in Japanese texts <i>Haruko Sanada</i>	183
Quantitative data processing in the ORD speech corpus of Russian everyday communication <i>Tatiana Sherstinova</i>	195
Complex investigation of texts with the system “StyleAnalyzer” <i>O.G. Shevelyov, V.V. Poddubnyj</i>	207
Retrieving collocational information from Japanese corpora: its methods and the notion of “circumcollocate” <i>Tadaharu Tanomura</i>	213
Diachrony of noun-phrases in specialized corpora <i>Nicolas Turenne</i>	223
Subject index	237
Author index	243
Authors’ addresses	247