# A Readability Checker
# Based on Deep Semantic Indicators

Tim vor der Brück and Sven Hartrumpf

Intelligent Information and Communication Systems (IICS)
FernUniversität in Hagen
58084 Hagen, Germany
{tim.vorderbrueck,sven.hartrumpf}@fernuni-hagen.de

**Abstract.** One major reason that readability checkers are still far away from judging the understandability of texts consists in the fact that no semantic information is used. Syntactic, lexical, or morphological information can only give limited access for estimating the cognitive difficulties for a human being to comprehend a text. In this paper however, we present a readability checker which uses semantic information in addition. This information is represented as semantic networks and is derived by a deep syntactico-semantic analysis. We investigate in which situations a semantic readability indicator can lead to superior results in comparison with ordinary surface indicators like sentence length. Finally, we compute the weights of our semantic indicators in the readability function based on the user ratings collected in an online evaluation.

**Keywords:** readability, understandability, semantics, indicator weights, linear optimization, linear regression

## 1   Introduction

Basically, a readability[1] checker has two major application areas. First, it can be used to automatically identify easy-to-read texts in a text corpus. In this case, it suffices to provide a global score which is usually calculated by a readability formula.

Second, a readability checker can be used to support authors to make their texts easy to read. In this case, more support is desirable than to compute only a global readability score. Instead, text passages which are difficult to read should be highlighted (e.g., the readability checker[2] of Rascu [1]). The calculation of a global score can here be useful too in order to give an estimation of the understandability of a text.

---

[1] In this paper, we use readability in the sense of understandability. We are aware that there exist other definitions where readability (or better: legibility) only relates to the form, but not to the contents of a text.

[2] Readability was not the only objective in this system. One further aspect was to ensure the fulfillment of certain formulation standards.

In this paper, we will discuss both application areas. Therefore we describe how semantic information can improve both the calculation of a global readability score and the identification of difficult text passages. Readability checkers can compute a global score by applying a readability formula on several indicator values.

## 2  Related Work

There are various methods to derive a numerical representation of text readability. One of the most popular readability formulas was created in 1948: the so-called Flesch Reading Ease [2]. The formula employs the average sentence length[3] and the average number of syllables per word for judging readability. The sentence length is intended to roughly approximate sentence complexity, while the number of syllables approximates word frequency since usually long words are less used. Later on, this formula was adjusted to German [3]. Despite of its age, the Flesch formula is still widely used.

Also, the revised Dale-Chall readability index [4] mainly depends on surface-type indicators. Actually, it is based on sentence length and the occurrences of words in a given list of words which are assumed to be difficult to read.

Recently, several more sophisticated approaches which use advanced NLP technology were developed. They determine for instance the embedding depth of clauses, the usage of active/passive voice or text cohesion [5–7]. The method of [8] goes a step beyond pure analysis and also creates suggestions for possible improvements.

As far as we know, all of those approaches are based on surface or syntactic structures but not on a truly semantic representation, like a semantic network as described here, which represents the cognitive difficulties for text understanding more adequately.

## 3  Semantic Networks

Semantic networks (SNs) of the MultiNet (Multilayered Extended Semantic Networks) formalism [9] allow to homogeneously represent the semantics of single words, phrases, sentences, texts, or text collections. Such SNs are chosen as the semantic representation in our DeLite readability checker described in this paper.

An SN node represents a concept, while an SN arc expresses a relation between two concepts. In MultiNet, each node is semantically classified by a *sort* from a hierarchy of 45 sorts. Furthermore, a node has an inner structure (depending on its sort) containing *layer features* like CARD (cardinality) and REFER (referential determinacy). Fig. 2 and Fig. 3 show the graphical form of SNs. They were generated by the WOCADI parser [10], which is employed in DeLite.

The WOCADI parser can construct SNs of the MultiNet formalism for German phrases, sentences, or texts. The text that is analyzed for readability is

---

[3] Throughout this work, sentence length is measured in words.

parsed sentence by sentence. During this process, SNs and syntactic dependency structures are built.

An important component of our deep syntactico-semantic analysis of natural language is HaGenLex, a semantically based computer lexicon [11]. This lexicon not only lists verb valencies, but also their syntactic and semantic types. Consider for example the German verb *essen* (*'eat'*). Sentences like *Die Birne isst den Apfel.* (*'The pear eats the apple.'*) are rejected because semantic selectional restrictions are violated. Besides this comprehensive lexicon with around 28,000 entries, we employ a flat lexicon, many name lexicons, and a sophisticated compound analysis to achieve the parser coverage required for applications like readability checkers.

Disambiguation is realized by specialized modules which work with symbolic rules and disambiguation statistics derived from annotated corpora. Currently, such modules exist for (intrasentential and intersentential) coreference resolution, the attachment of prepositional phrases, and the interpretation of prepositional phrases.

## 4 Conception of Our Readability Checker

Readability can be measured by way of numerous *readability criteria*. Each criterion (like semantic complexity) can be realized or approximated by one or more operable (i.e., implementable) *readability indicators* (like *Number of propositions per sentence*, *Maximum path length in the SN*, etc.). Note that an indicator can only be applied on a specific type of text segments which we call the *segment type*[4] of this indicator, e.g., the indicator *Number of propositions per sentence* can only be applied on an entire sentence, but not on single words. We differentiate between the segment types *word*, *phrase*, *sentence*, and *text*.

### 4.1 Calculating a Global Readability Score

In DeLite the calculation of the global readability score is done in several steps (see Fig. 1):

- Segmentation: In the first step, the entire document is segmented into words, phrases, and sentences based on the parser results.
- (Basic) Calculation: Indicator values are calculated for each segment the indicators are associated to, e.g., the indicator *Number of concepts in a compound* calculates one value for every word, the indicator *Sentence length* for every sentence.
- Aggregation: For each indicator, its values associated to text segments are averaged. This average is called the aggregated indicator value.
- Normalization: To combine indicators of different types their values have to be transformed to a common value range. In DeLite, the aggregated indicator values are all mapped to the interval from zero to one.

---

[4] In some rare cases the applicability is further restricted, e.g., the indicator *Number of reference candidates* is not applicable to all kinds of words, but only to pronouns.
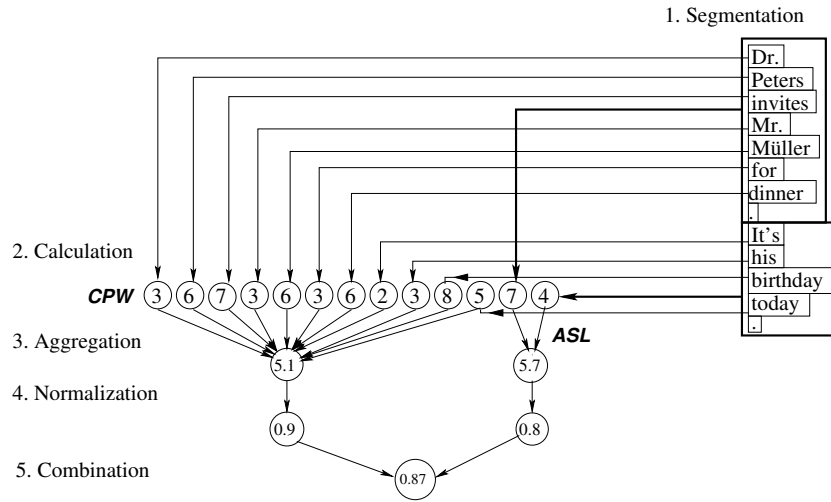
**Fig. 1.** Calculation of a global readability score with the two indicators *Number of characters per word (CPW)* and *Average sentence length (ASL)*.

- Combination: In the last step, a global readability score is determined by calculating a weighted sum of all aggregated and normalized indicator values. All weights are non-negative and sum up to one.

### 4.2  Highlighting Text Segments

We compute an indicator value for each text segment this indicator is applicable to. If that value exceeds a certain threshold, the associated text segment is highlighted. For example, if the threshold for the indicator *Number of concept nodes in the SN* is 10, all sentences having 11 or more concept nodes will be highlighted.

We experienced that this approach did not suffice. Sometimes it is important for an exact understanding of the readability problem to highlight additional text segments. We call these text segments *supplementary highlight segments* in contrast to the *primary highlight segments* which directly refer to the found readability defect. (See Sect. 5.3 for an example.) Note that the segment type of a supplementary highlight segment does not have to match the segment type of the associated primary highlight segment.

In the following section, we describe some of the most important semantic readability indicators. For more motivation and references to the literature (e.g., from psycholinguistics) please see [12] and [13].

# 5 Semantically Oriented Readability Indicators

## 5.1 Abstract and Concrete Nouns

A high proportion of abstract nouns can deteriorate text readability [14]. A noun is considered as abstract if it does not directly refer to a visible object. The binary information whether a noun is abstract or not is available from our semantically oriented lexicon. The annotation is made on concepts and not on words since a word can have both abstract and concrete readings. For example, the German word *Platz* can mean *a place in a city* (like a plaza) which is a visible, concrete object. Alternatively, it can mean *space* like in the sentence: *Im Englisch-Kurs ist kein Platz frei.* (*'There is no space left in the English course.'*).

## 5.2 Negation

Negations can make a sentence more difficult to understand [14] and should be avoided if a positive formulation is possible. There exist many possibilities to convey negation in German [15]. Negation can be expressed by special words, e.g., *nicht* (*'not'*) and *niemals* (*'never'*), or prefixes, e.g., *unmöglich* (*'impossible'*) is the antonym of *möglich* (*'possible'*). While special words are quite easy to recognize, this is not the case for negation prefixes. First, such a prefix is not trivial to recognize, e.g., the German word *unterirdisch* does not contain the negation prefix *un*, but the prefix *unter* (*'under'*), which has a completely different meaning. Second, in some cases a word contains actually a negation prefix, but it is not used as a negation, e.g., the adjective *unheimlich* (*'weird'*) is not an antonym of *heimlich* (*'secret'*). However, if semantic information is available, this problem can be handled quite easily. Consider we have some word $w$ which is the concatenation of the prefix *un* and a word $v$. We can infer that $w$ is a negated adjective, if $w$ is an antonym of $v$ (which means that the lexicon contains an ANTO (antonymy) relation connecting $v$ and $w$). Note that there exist several algorithms to extract semantic relations like ANTO by analyzing large text corpora. These methods would save the work to manually add ANTO relations to the lexicon; however, for cases like *unheimlich* and *heimlich* above, special treatment (or manual correction) is needed.

A special case of negations are *double negations*. A sentence contains a double negation if a similar (but not the same) semantics can be achieved by dropping two negations occurring in this sentence. This effect takes place if one negation is in the scope of another. Note that there are also sentences which contain triple or quadruple negations, e.g., the sentence *Ich glaube nicht, dass Peter nicht denkt, dass der Film nicht uninteressant ist.* (*'I do not believe that Peter does not think that the movie is not uninteresting.'*) contains a quadruple negation. In almost all cases, double negations are redundant and should be avoided. A double negation can relate to a sentence, to a phrase, or only to a word. Our readability checker can recognize several different kinds of double negations, e.g., a double negation occurs in a sentence if the sentence node is associated to the facticity (layer feature FACT) *nonreal* and is connected to the modality *non.0* by a MODL (modality) relation; see [9] for details on the semantic representation.

### 5.3  Indicators Concerning Anaphors

Several readability problems can concern anaphors. Consider the sentence: *Dr. Peters lädt Herrn Müller zum Essen ein, da heute sein Geburtstag ist.* (*'Dr. Peters invites Mr. Müller for dinner since it is his birthday today.'*). The possessive determiner *sein* (*'his'*) can either relate to the antecedent candidate *Dr. Peters* or to the antecedent candidate *Mr. Müller*. For a better understanding this sentence should be reformulated, e.g., by repeating either *Dr. Peters* or *Mr. Müller*. Thus we introduced a readability indicator counting the number of possible antecedents for each anaphor. In DeLite, the anaphor is marked as primary and the possible antecedents as supplementary highlight segments if this indicator value exceeds the associated threshold (e.g., 1).

Furthermore, an anaphoric reference can be difficult to resolve if the antecedent is too far away from the anaphor. The distance can be measured in words, sentences, or—more semantically and psycholinguistically motivated— by intervening entities (or discourse referents). Finally, we also use an indicator to check if there exists at least one antecedent for each anaphor.

### 5.4  Number of Propositions per Sentence

A further measure for sentence complexity is the number of SN nodes which bear the semantic sort *si* (situation, like *to discuss*) or *abs* (abstract situation, for nominalized verbs like *discussion*) or one of their subsorts. Such nodes correspond to the propositions in a given sentence. This indicator is correlated to the sentence length since a long sentence usually contains also several propositions. However, this is not always the case. Consider for example the following long sentence: *Anwesend waren Dr. Schulz, Dr. Peters, Herr Werner, Frau Brand, Herr Mustermann, Herr Frank, Dr. Grainer, [. . . ].* (*'Dr. Schulz, Dr. Peters, Mr. Werner, Mrs. Brand, Mr. Mustermann, Mr. Frank, Dr. Grainer, [. . . ] were present.'*) which contains only a single proposition. Long item lists usually do not degrade readability [16]. Therefore in such situations the readability can more appropriately be judged by the indicator *Number of propositions per sentence* than by *Average sentence length*.

Also the opposite effect can be found: a quite short sentence can contain many propositions (for example expressed by participle constructions). The indicator *Average sentence length* would not be violated, while the sentence is definitely hard to read, e.g., *The man* running *downhill and* meeting *the colleague* walking *to the office* fell *over a dog* chased *by a boy*. This sentence contains five propositions and is definitely hard to understand.

### 5.5  Maximum Path Length

We measure the length of the longest path that the SN contains which is based on the assumption that information is often more difficult to understand if the constituents depend on each other and therefore a sequential interpretation is necessary. Consider for example the easy-to-read sentence *Ich besuche*
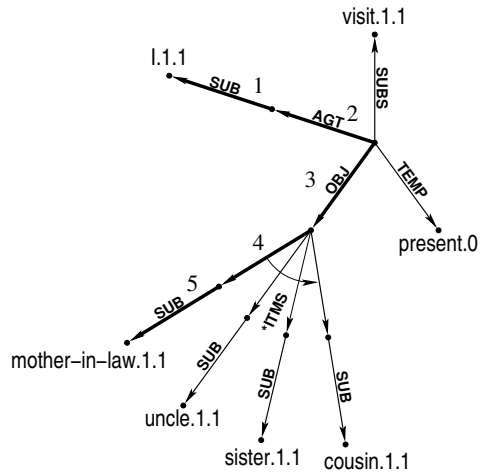
**Fig. 2.** Semantic network for the sentence *Ich besuche die Schwiegermutter, den Onkel, die Schwester und die Cousine.* (*'I visit the mother-in-law, the uncle, the sister, and the cousin.'*). One longest path, not taking into account the direction of arcs, is printed in bold face.

*die Schwiegermutter, den Onkel, die Schwester und die Cousine.* (*'I visit the mother-in-law, the uncle, the sister, and the cousin.'*). Since the constituents in the coordination do not depend on each other they can be interpreted in parallel which makes the sentence easy to understand. The length of the longest path in the SN is 5 which is still rather short (see Fig. 2). However, this is not the case for the following sentence where the constituents have to be interpreted sequentially: *Ich besuche die Schwiegermutter des Onkels der Schwester der Cousine.* (*'I visit the mother-in-law of the uncle of the cousin's sister.'*). Similar effects can be observed in connection with negations where the special phenomena of double negations can emerge (see Sect. 5.2). For this sentence, the length of the longest path is 7 (see Fig. 3). Thus, sequentially interpreted sentences usually lead to longer paths in the SN.

### 5.6 Semantic Network Quality

The case that the SN for some sentence could not be constructed or is assigned a low quality score is often caused by the fact that the associated sentence is syntactically or semantically complex or even incorrect. Thus we provide an indicator for this information. Note that this indicator is not purely semantic since the construction of the SN can fail if the syntactic structure of the sentence is invalid.
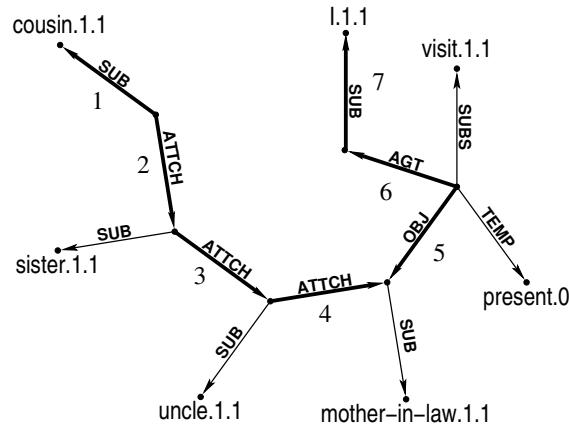
**Fig. 3.** Semantic network for the sentence *Ich besuche die Schwiegermutter des Onkels der Schwester der Cousine.* (*'I visit the mother-in-law of the uncle of the cousin's sister.'*). The longest path, not taking into account the direction of arcs, is printed in bold face.

### 5.7 Passive Construction

The syntactic indicator *Passive* was enriched with semantic information leading to the new indicator *Passive with agent.* Usually sentence formulations in active voice are easier to understand than equivalent formulations in passive voice [14]. To convert a sentence into active voice the direct object and the subject have to change roles. We call the new subject the *semantic subject.* Passive constructions are very common in German. Thus we want to highlight a passive sentence (or reduce the readability score) only if it is obvious that an active formulation would be better.

There exist some exceptions to the rule that active formulations should be preferred. In some cases the semantic subject might not be known (or might be irrelevant), e.g., *Peter wurde rechtzeitig benachrichtigt.* (*'Peter was informed on time.'*). In this case, the impersonal pronoun *man* (*'one'*) can be inserted to convert the sentence into active: *Man benachrichtigte Peter rechtzeitig.* However, this formulation is usually not better than the original. Moreover, sometimes a passive formulation will be preferred if the semantic subject is neither a human being nor an animal. For example, the sentence *Peter wurde vom Blitz erschlagen.* (*'Peter was struck by a lightning.'*) need not be converted into *Der Blitz erschlug Peter.* (*'The lightning struck Peter.'*).

Since a complete linguistic treatment of all cases is not trivial we used a heuristic. We only penalized passive if the semantic subject is uttered and is connected to the sentence by the semantic relation AGT (agent). In this case, the semantic subject usually performs some sort of action and an active formulation should always be possible. This heuristic conforms to [17] who propose that an active formulation should be preferred if the sentence is agent-oriented.
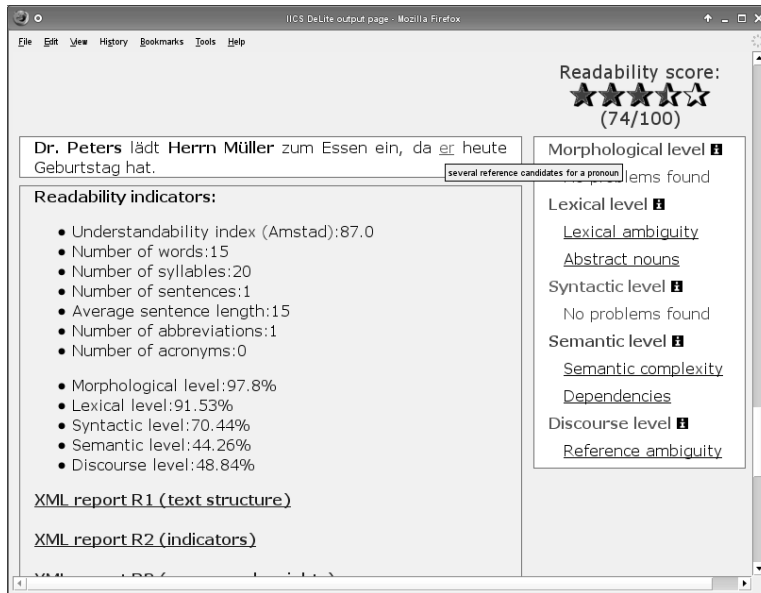
**Fig. 4.** Screenshot of DeLite's GUI in which a referential pronoun ambiguity is indicated for the sentence *Dr. Peters lädt Herrn Müller zum Essen ein, da er heute Geburtstag hat.* ('*Dr. Peters invites Mr. Müller for dinner, since it is his birthay today.*'; literally: '*. . . , since he has birthday today.*').

### 5.8 Other Semantic Readability Indicators

We evaluated further semantic indicators. For instance, the inverse concept frequency is determined which is based on readings (as determined by word sense disambiguation in the parser) instead of word forms. This indicator can detect cases where a reading is rare but the word (as a whole) is not rare.

We also introduced an indicator determining the average number of arcs the discourse entities of the SN were connected to (*Connectivity of discourse entities*), where the discourse entities are identified as SN nodes with the ontological sort *object* [9, p. 409–411]. For concessive and causal clauses, DeLite counts the causal and concessive relations in a chain.

## 6 User Interface

We provided a graphical user interface (GUI) for the readability checker DeLite [12] which displays a global readability score and highlights (by color) text passages that are difficult to read according to at least one indicator (see Fig. 4). If the user moves the mouse on such a text passage, the readability problem type will be described briefly. Supplementary highlight segments (if any) are printed in bold face if the user clicks on the colored text passage. In the upper right

corner, a global readability score is provided which is calculated by a readability formula over all readability indicators.

## 7 Evaluation

We evaluated our algorithm as implemented in DeLite on a text corpus of 500 texts from the local administration domain. 315 users participated in the readability study, 43.1 % of them were female and 56.9 % male.

Almost 70 % of the participants were between 20 and 40 years old; the number of participants over 60 was very small (circa 3 %). The participants were mainly well-educated. 58 % of them owned a university or college degree. There is none who had no school graduation at all. The participants of the evaluation belonged to a large variety of professions, e.g., software developers, scientists, physicians, linguists, pharmacists, administrators, psychologists, and musicians. Each participant rated the readability of several texts on a 7-point Likert scale [18].

We determined the weight of each indicator in our readability formula using both linear optimization [19] and linear regression with the Lagrange restriction [20] that indicator weights sum up to one (see Table 1). Both methods represent the readability score, as determined by the participants of the readability study, as a weighted sum of normalized indicators and estimated the weights in such a way that the mean absolute error (linear regression: mean squared error) is minimized [21]. Only 13 indicators (of 53 indicators) were assigned a weight greater than zero. The evaluation showed that deep semantic and syntactic indicators have quite comparable weights to traditional surface type indicators. The weights of the semantic indicators are expected to further improve if parser quality and coverage increases. Note that the weights should be seen with caution since changes in the parser can have a serious impact on them.

In a second experiment, we replaced the indicator *Abstract noun*, which was assigned a weight of zero, by a combination of the indicators *Abstract noun* and *Inverse concept frequency* in such a way that the indicator value is assigned to zero if the noun is concrete and to the inverse concept frequency otherwise. Afterwards, the weights were redetermined using the above-mentioned optimization algorithms. In this second experiment, the combined indicator was assigned the weight 9.7 % (7.8 % for linear optimization ) while the weight of the now strongly correlated indicator *Inverse concept frequency* did not decrease.

Semantic indicators with the strongest correlation to the user ratings were *Semantic network quality*, *Inverse concept frequency*, *Maximum path length in the SN*, *Connectivity of discourse entities*, *Number of propositions per sentence*, and several anaphora related indicators. We noticed, however, that indicator correlation was not very reliable for estimating indicator importance in the readability function since, on the one hand, quite strongly correlated indicators can have a low weight if they are highly correlated to other indicators. On the other hand, indicators with rather weak correlation can have a considerable impact in the readability function if they are only weakly correlated to the other indicators.

**Table 1.** Selected indicator weights; Sur=surface type indicator, Syn=syntactic indicator, Sem=semantic indicator.

| Indicator | Type | Weight (%) | |
| --- | --- | --- | --- |
| | | Lin. regression | Lin. optimization |
| Average sentence length | Sur | 34.1 | 35.0 |
| Semantic network quality | Sem/Syn | 20.3 | 27.6 |
| Number of syllables per word | Sur | 12.6 | 10.8 |
| Number of words per NP | Syn | 6.3 | 3.1 |
| Inverse concept frequency | Sem | 6.0 | 6.0 |
| Word form frequency | Sur | 5.9 | 1.4 |
| Maximum path length in the SN | Sem | 4.7 | 2.9 |
| Conditional relations in a chain | Sem | 1.3 | 2.1 |
| Distance verb complement | Syn | 0.9 | 0.8 |
| Reference distance of a pronoun in words | Sem | 0 | 1.2 |
| *3 other indicators* | *All* | 7.9 | 9.1 |
| Number of characters per word | Sur | 0 | 0 |
| *39 other indicators* | *All* | 0 | 0 |

Finally, the DeLite readability index was compared to a baseline: the Amstad readability index [3]. Applied on our test corpus this readability index reached a correlation with the user ratings of 0.187 which is far below the DeLite correlation of 0.509. However, this difference is mainly caused by the fact that the parameters of the Amstad readability index were derived by analyzing newspaper texts, which differ considerably from documents of local administration used here. Thus, we additionally determined a readability index resulting from employing a linear optimization only on the two indicators of the Amstad readability index, i.e., *Average sentence length* and *Number of syllables per word*. The correlation increased considerably to 0.458 but is still clearly outperformed by the DeLite index.

## 8 Conclusion and Future Work

We proposed a new kind of readability indicators which are semantic and predominantly operate directly on semantic representations (SNs). We further investigated indicator weights and correlations of indicators and user ratings. The evaluation showed that, although the SN could not be constructed for several sentences of our domain-specific corpus, semantic indicators can often yield scores that are more accurate than traditional, surface-oriented readability indicators. Therefore we expect that semantic readability indicators will play an important role for future readability checkers.

## Acknowledgments

## References

1. Rascu, E.: A controlled language approach to text optimization in technical documentation. In: Proceedings of KONVENS 2006, Konstanz, Germany (2006) 107–114
2. Flesch, R.: A new readability yardstick. Journal of Applied Psychology **32** (1948) 221–233
3. Amstad, T.: Wie verständlich sind unsere Zeitungen? PhD thesis, Universität Zürich, Zürich, Switzerland (1978)
4. Chall, J., Dale, E.: Readability Revisited: The New Dale-Chall Readability Formula. Brookline Books, Brookline, Massachusetts (1995)
5. McCarthy, P., Lightman, E., Dufty, D., McNamara, D.: Using Coh-Metrix to assess distributions of cohesion and difficulty: An investigation of the structure of high-school textbooks. In: Proc. of the Annual Meeting of the Cognitive Science Society, Vancouver, Canada (2006)
6. Heilman, M.J., Collins-Thompson, K., Callan, J., Eskenazi, M.: Combining lexical and grammatical features to improve readability measures for first and second language texts. In: Proceedings of the Human Language Technology Conference, Rochester, New York (2007)
7. Segler, T.M.: Investigating the Selection of Example Sentences for Unknown Target Words in ICALL Reading Texts for L2 German. PhD thesis, School of Informatics, University of Edinburgh (2007)
8. Chandrasekar, R., Srinivas, B.: Automatic induction of rules for text simplification. Technical Report IRCS Report 96-30, University of Pennsylvania, Philadelphia, Pennsylvania (1996)
9. Helbig, H.: Knowledge Representation and the Semantics of Natural Language. Springer, Berlin, Germany (2006)
10. Hartrumpf, S.: Hybrid Disambiguation in Natural Language Analysis. Der Andere Verlag, Osnabrück, Germany (2003)
11. Hartrumpf, S., Helbig, H., Osswald, R.: The semantically based computer lexicon HaGenLex – Structure and technological environment. Traitement automatique des langues **44**(2) (2003) 81–105
12. Hartrumpf, S., Helbig, H., Leveling, J., Osswald, R.: An architecture for controlling simple language in web pages. eMinds: International Journal on Human-Computer Interaction **1**(2) (2006) 93–112
13. Jenge, C., Hartrumpf, S., Helbig, H., Nordbrock, G., Gappa, H.: Description of syntactic-semantic phenomena which can be automatically controlled by NLP techniques if set as criteria by certain guidelines. EU-Deliverable 6.1, FernUniversität in Hagen (2005)
14. Groeben, N.: Leserpsychologie: Textverständnis – Textverständlichkeit. Aschendorff, Münster, Germany (1982)

15. Drosdowski, G.: Duden - Grammatik der deutschen Gegenwartssprache. Dudenverlag, Mannheim, Germany (1995)
16. Langer, I., von Thun, F.S., Tausch, R.: Sich verständlich ausdrücken. Reinhardt, München, Germany (1981)
17. Helbig, G., Kempter, F.: Das Passiv. Zur Theorie und Praxis des Deutschunterrichts für Ausländer. Langenscheidt, Berlin, Germany (1997)
18. Likert, R.: A technique for the measurement of attitudes. Archives of Psychology **140** (1932) 1–55
19. Bertsimas, D., Tsitsiklis, J.: Introduction to Linear Optimization. Athena Scientific, Belmont, Massachusetts (1997)
20. Greene, W.: Econometric Analysis. Prentice Hall, Englewood Cliffs, New York, USA (1993)
21. vor der Brück, T., Leveling, J.: Parameter learning for a readability checking tool. In Hinneburg, A., ed.: Proceedings of the LWA 2007 (Lernen-Wissen-Adaption), Workshop KDML. Gesellschaft für Informatik, Halle/Saale, Germany (2007)