

# Understanding Document Understandability Through Correlation Analysis

Anonymous

No Institute Given

## 1 Introduction

Recently, a number of works rely on the output of readability formulas to deem if a Web page is easy or hard to understand (e.g., [7, 6, 17, 14, 1, 11]). However, as discussed in Chapter ??, preprocessing steps are crucial for the interpretation of readability formula results. When applying them to assess the understandability of Web document, the same readability formula may yield results that vary from *suitable even for kids* to *understandable only if you are an experts* with only minor modifications in the preprocessing steps. In this chapter, we investigate if any of these two interpretations is correct analysing the correlations of various understandability estimators and human assessments. For that, we take advantage of the understandability assessments made in recent CLEF eHealth campaigns ([13, 21]).

During the CLEF eHealth 2015 and 2016 campaigns [13, 21], in addition to topical relevance, medical students and professional serving as assessors were instructed to assess how easy to understand and how much they would trust the information contained in each assessed document. The assessment tool used in these campaigns is shown in Figure ?? to enhance the reader understanding.

The understandability assessments for the topically relevant documents are used throughout this chapter. We limit our analysis to the relevant documents to reduce the noise associated with understandability assessments, as the use of ClueWeb-12 B in 2016 campaign resulted in documents from a broad range of topics being retrieved, not only medical ones, and we would like to keep in our analysis only documents related to medicine.

As also depicted in Figure ??, understandability assessments were given on a 4-label scale (scores from 0 to 3) in 2015 and on a 0-100 scale in 2016. Easy to understand documents had assessments closer to 0 while hard to understand document had their assessments close to the maximal value, 3 in 2015 and 100 in 2016. Altogether there are 1,452 documents assessed for understandability for CLEF eHealth 2015 and 3,320 for 2016 (see a detailed analysis in Appendix ??, for example, Figures ?? and ?? report the assessment distribution). We use the Pearson, Kendall Tau and Spearman's Rank correlation in this work. The first one is used to calculate the strength of two linearly related variables, while the last two are rank correlation, i.e., they act on the rank of the variables instead of their values. We opt to report on all three variables as all three are equally used in the understandability literature, and thus we allow that other researchers can compare our results across.

Add reason for using all 3 correlations

We start our experiments in this chapter analysing the correlation of each readability formula to the human assessments, our ground truth, in Section 2. We then study other understandability estimators going far beyond the readability formulas. Those are presented in Section 3. Similarly, to our experiments with readability formulas, we empirically investigate in Section 4 the patterns that lead to effective understandability estimators. In Section 5 we compare the effect of different preprocessing methods in the light of more than 100 understandability estimators. Our findings and highlights are summarized in Section 7.

## 2 Which Readability Formula To Use

Our previous experiments summarized in Figure ?? found that the use of Naive preprocessing was associated with larger variances in the understandability score estimated by readability formulas. We plot in Figures 1 and 2 the correlation scores of each traditional readability metric with the human assessments made in CLEF 2015 and 2016, respectively. We observe that the Naive preprocessing also results in the lowest correlation, no matter which correlation measure or readability formula is used. Also, when the Naive preprocessing is used, the variant DoNotForcePeriod yields higher correlations than the variant ForcePeriod, but when using a higher quality HTML cleaner, such as Juxtext or Boilerplate, the results indicate that the use of ForcePeriod should be preferred.

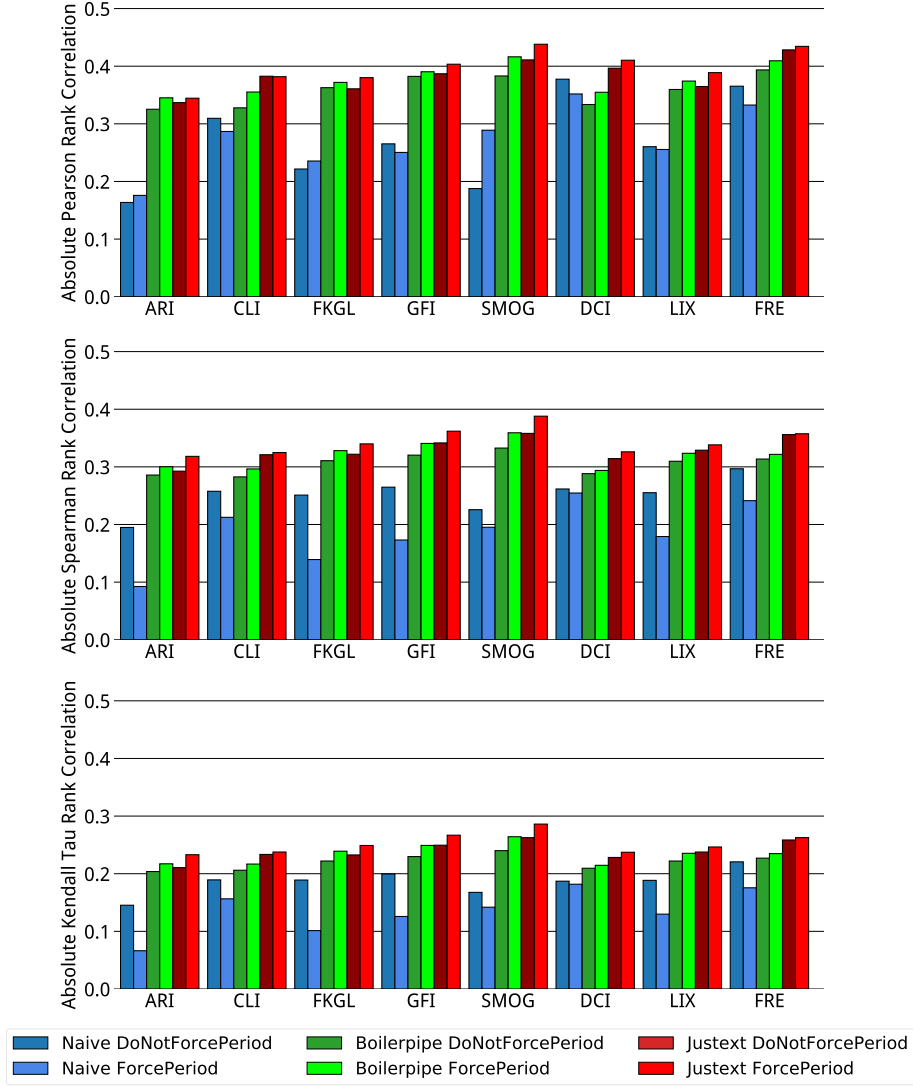
Among all the readability formulas and preprocessing methods, SMOG with ForcePeriod preprocessing and Dale-Chall Index with DoNotForcePeriod are the best ones respectively for 2015 and 2016. Although there is no single best readability measure or best preprocessing strategy in all scenarios, CLI and FRE with Juxtext are stable options, with correlation as high as the best ones in both campaigns. Thus, we advice for the use of CLI, as it has also been shown to be the most robust measure to variances due to use of ForcePeriod or DoNotForcePeriod in Figure ??.

## 3 (More) Understandability Estimators

The correlation of readability formulas as shown in Figures 1 and 2 is not strong, with no correlation being higher than 0.5. One of our intents with the rest of this chapter is comparing the correlation of the traditional readability formulas with other methods for understandability estimation, including an evaluation of other humans performing the same task. For that, we devise and group several methods into semantically related groups which will be following presented. We summarize all methods in Table 2.

### 3.1 Traditional Readability Formulas

This group contains all the readability formulas listed in Chapter ?? (Table ??) and additional readability formulas that were excluded for not being on the same value scale as the other ones.

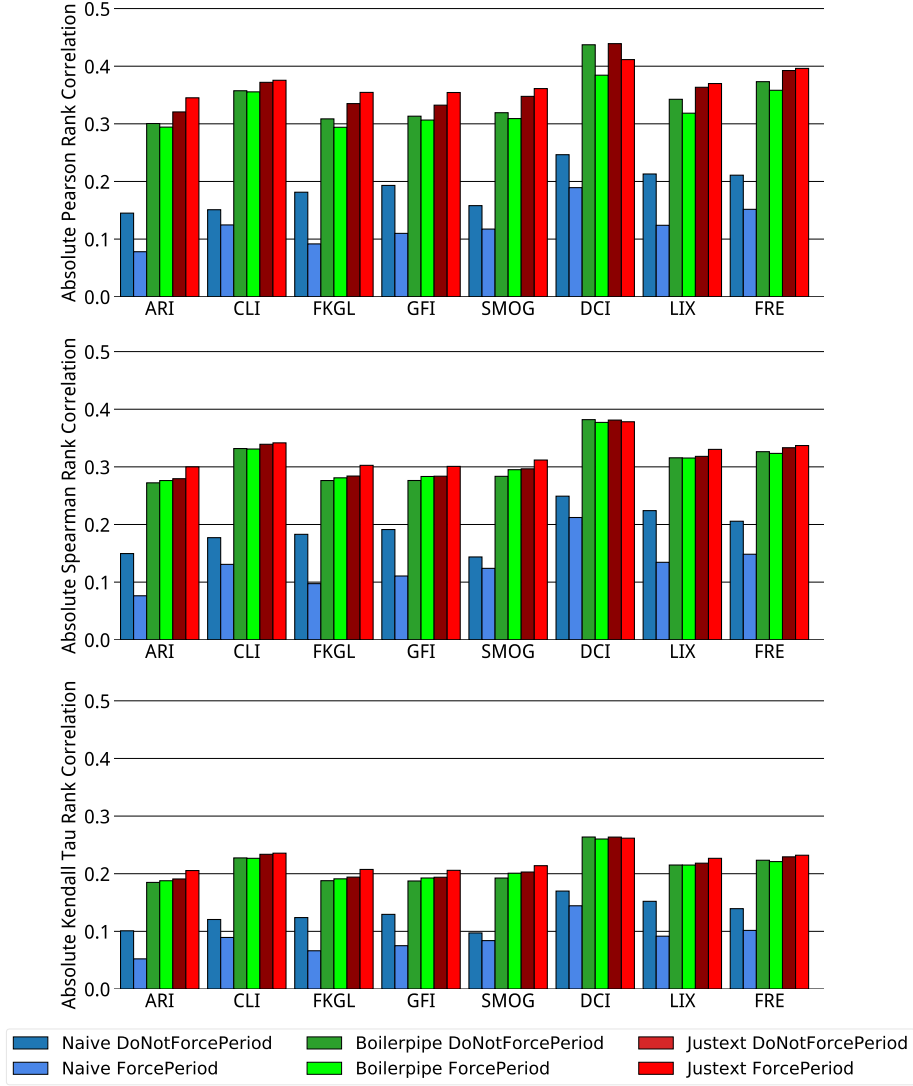


**Fig. 1.** Correlation of different readability measures and the understandability scores collected in CLEF eHealth 2015

### 3.2 Raw Components of Readability Formulas

This group comprises the building blocks that make up the traditional readability measures. Some examples include the average number of characters per word or the average number of syllables in a sentence<sup>1</sup>.

<sup>1</sup> Words were divided into syllables using the Python package Pyphen <http://pyphen.org/>



**Fig. 2.** Correlation of different readability measures and the understandability scores collected in CLEF eHealth 2016.

### 3.3 General Medical Vocabularies

This group includes methods such as the number of words with a medical prefix or suffix, i.e. beginning or ending with Latin or Greek particles (e.g., amni, angi, algia, arteri), acronyms<sup>2</sup> or medical vocabularies such as the International

<sup>2</sup> The acronym list was obtained from the ADAM database [20]

Statistical Classification of Diseases and Related Health Problems (ICD), Drugbank and the OpenMedSpel dictionary<sup>3</sup>. Methods listed here were matched with documents using a simple keywords matching.

### 3.4 Consumer Vocabulary Features

the Consumer Health Vocabulary (CHV) is a prominent medical vocabulary dedicated to mapping consumer (layperson) vocabulary to technical terms. It attributes a score for each of its concepts with respect to their difficulty, with lower/higher scores for harder/easier concepts. We used MetaMap once again to map the content of Web documents, as done in Chapters ?? and ?. We further use MetaMap options to also filter only concepts identified as symptoms or diseases, using the same definitions from Section ??.

### 3.5 Expert Vocabulary Features

The hierarchy of Medical Subject Headers (MeSH) was previously used in the literature to identify hard concepts, assuming that a concept that is deep in the hierarchy is harder than a shallow one [18]. As done with CHV, we used MetaMap to map the content of Web documents to MeSH and explore symptoms and disease concepts separately.

### 3.6 Natural Language

This group comprises commonly used metrics in the natural language processing field: the ratio of part-of-speech (POS) classes, the number of entities in a text, the sentiment polarity and the ratio of words found in English vocabularies. The Python package NLTK 3.2<sup>4</sup> was employed for sentiment analysis and POS tagging. The GNU Aspell<sup>5</sup> dictionary was used as a standard English vocabulary and a stopword list was built by merging the stopword lists of the Indri<sup>6</sup> and Terrier<sup>7</sup> toolkits.

### 3.7 HTML Features

The aim of this group is to represent a web page by its HTML content. We hypothesize that a Web page rich of images or with its content well summarized in tables can potentially ease hard subjects such as medicine. We identify a large number of HTML tags in this group with the Python library BeautifulSoup v4.4<sup>8</sup>.

<sup>3</sup> <http://extensions.openoffice.org/en/project/openmedspel-en-us>

<sup>4</sup> <http://www.nltk.org/>

<sup>5</sup> <http://www.aspell.net/>

<sup>6</sup> <http://www.lemurproject.org/indri/>

<sup>7</sup> <http://www.terrier.org/>

<sup>8</sup> <http://www.crummy.com/software/BeautifulSoup/>

### 3.8 Word Frequency Features

Common and known words are usually frequent words, while unknown and obscure words are rare. This idea is implemented in readability formulas such as the Dale-Chall index which uses a list of common words and counts the number of words that fall outside this list [5]. In this work we model word frequency in a straightforward manner: we sort the frequency of all words in a corpus and normalize the ranking of word frequency such that values close to 100 are attributed to common words and values close to 0 to rare words. We explore three different corpora in this work:

- Medical Reddit: Reddit is an Internet forum with a sizable user community which is responsible for generating content. Any user can start a discussion receiving replies from any other user. This discussion forum is intensively used for health purposes, for example in the Reddit community AskDocs licensed nurses and doctors (subject to user identity verification) advise help seekers free of charge. We selected six of such communities (medical, AskDocs, AskDoctorSmeeee, Health, WomensHealth, Mens\_Health) and downloaded all user interactions using the Python Reddit API Wrapper (PRAW<sup>9</sup>), v5.1. In total 43,018 discussions were collected.
- PubMed Central: PubMed Central (PMC) is an online digital database of freely available full-text biomedical literature playing a similar role to physicians as the ACM Digital Library does to computer scientists. We used in this work the same collection crafted for the TREC Clinical Decision Support Track 2014 and 2015 (TREC-CDS)<sup>10</sup> consisting of 733,138 articles.
- Medical English Wikipedia: we filtered articles from a Wikipedia dump<sup>11</sup> (May 1st 2017), that contained an Infobox<sup>12</sup> in which at least one of the following words appeared as a property: ICD10, ICD9, DiseasesDB, MeSH, MeSHID, MeshName, MeshNumber, GeneReviewsName, Orphanet, eMedicine, MedlinePlus, drug\_name, Drugs.com, DailyMedID, LOINC. Figure 3 illustrates a Wikipedia page that is marked as medical because of its Infobox entries. This idea was successfully implemented in Soldaini et al. [16] and our filtering process resulted in a collection of 11,942 articles. Note that this procedure highly favors precision over recall.

A summary of the statistics of these three collections is reported in Table 1. In order to calculate word frequency, we removed words that occur less than 5 times in a corpus. Finally, we ignore out of vocabulary (OV) words in our calculations, unless it is explicitly stated.

<sup>9</sup> <https://praw.readthedocs.io/>

<sup>10</sup> <http://www.trec-cds.org/>

<sup>11</sup> <https://dumps.wikimedia.org/enwiki/>

<sup>12</sup> A Wikipedia infobox is a template containing structured information that appear on the right of Wikipedia pages to summarize key aspects of concepts

Article
Talk
Read
Edit
View history
Search Wikipedia

## Hyperthermia

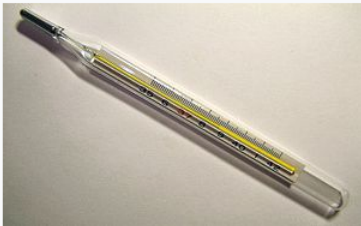
From Wikipedia, the free encyclopedia

**Hyperthermia** is elevated body temperature due to failed [thermoregulation](#) that occurs when a body produces or absorbs more [heat](#) than it dissipates. Extreme temperature elevation then becomes a [medical emergency](#) requiring immediate treatment to prevent disability or death.

The most common causes include [heat stroke](#) and adverse reactions to drugs. The former is an [acute temperature elevation](#) caused by exposure to excessive heat, or combination of heat and humidity, that overwhelms the heat-regulating mechanisms. The latter is a relatively rare side effect of many drugs, particularly those that affect the [central nervous system](#). [Malignant hyperthermia](#) is a rare complication of some types of [general anesthesia](#).

Hyperthermia differs from [fever](#) in that the body's [temperature set point](#) remains unchanged. The opposite is [hypothermia](#), which occurs when the temperature drops below that required to maintain normal metabolism. The term is from [Greek](#) ὑπέρ, *hyper*, meaning "above" or "over", and θερμός, *thermos*, meaning "hot".

Hyperthermia



An analog [medical thermometer](#) showing a temperature of 38.7 °C (101.7 °F)

**Classification and external resources**

<b>Specialty</b>	Critical care medicine
<b>ICD-9-CM</b>	780.6 <a href="#">↗</a>
<b>DiseasesDB</b>	18924 <a href="#">↗</a>
<b>MeSH</b>	D005334 <a href="#">↗</a>

[edit on Wikidata]

**Fig. 3.** Wikipedia page on hyperthermia. A rectangular red box identify the Infobox on the right hand side containing entries for Specialty, ICD-9-CM, DiseasesDB and MeSH.

### 3.9 Machine Learning on Text - Regressors and Classifiers

In a recent survey, Kevin Collins-Thompson reports that the future of understandability estimation relies on Machine Learning [4]. A challenge in using Machine Learning in this task is defining the background corpora used as training set. A possible setup for our work could have used CLEF 2015 assessments to learn a model for CLEF 2016 and vice-versa, but instead, we opt for a more reusable solution for the medical/health domain. We employed the three datasets explained in Section 3.8 and assume different labels according to the average difficulty of documents in these collections:

- Medical Reddit (label 1): Documents in this collection are expected to be written in a colloquial style, and thus the easiest to understand. All the conversations are in fact explicitly directed to assist inexperienced health consumers;
- Medical English Wikipedia (label 2): Documents in this collection are expected to be less formal than scientific articles, but more formal than a Web forum like Reddit;
- PubMed Central (label 3): Documents from this collection are expected to be written in a highly formal style, as the target audience of these documents are physicians, nurses and researchers in the biomedical domain.

**Table 1.** General statistics for the auxiliar collections used in this work

Statistic	Medical Wikipedia	Medical Reddit	PubMed Central
Number of Docs.	11,868	43,019	733,191
Number of Words	10,655,572	11,978,447	144,024,976
Number of Unique Words	467,650	317,106	2,933,167
Avg. Words per Doc.	898.90 $\pm$ 1351.76	278.45 $\pm$ 359.70	227.22 $\pm$ 270.44
Avg. Char per Doc.	5107.81 $\pm$ 7618.57	1258.44 $\pm$ 1659.96	1309.11 $\pm$ 1447.31
Avg. Char per Word	5.68 $\pm$ 3.75	4.52 $\pm$ 3.52	5.76 $\pm$ 3.51

Based on word counts in documents in these three collections, we model two different tasks a classification one and a regression task. Different labels for the regression could be employed, for example, a label 5 to PubMed Central documents to emphasize that these documents are explicitly made for expert users. We did not explore the effects of different labels in this work, it is left as future work.

#### 4 Top Measures from Each Group

We correlated each individual understandability estimator listed in Table 2 with the human assessments collected in CLEF eHealth 2015 and 2016 campaigns. We report in Table 3 the best metric for each group according to Pearson, Spearman or Kendall correlation. For some groups, such as the readability formula group, the highest correlated metric was the same for different correlation measure: SMOG Index in CLEF eHealth 2015 and Dale-Chall Index in 2016. We highlight the top score value of each correlation measure in each group. Note that there is no single case in which three different metrics were the top correlated for each different correlation measure.

hypotesis that kendaatl tau and spearman always point to the same winner

Interestingly, Table 3 shows that the polysyllable words, best formula component metric for CLEF 2015 data, is the main metric for the SMOG formula, the best readability formula for CLEF 2015. Likewise, the number of difficult words, best formula component metric for CLEF 2016, is the main metric for Dale-Chall index, the best readability formula for CLEF 2016.

The top correlation for MeSH group, number of MeSH concepts, reaches much lower correlation than the top correlation metric for the CHV group, the scores of CHV concepts. The dominating metrics for the Natural Language group are the number of pronouns, the number of stopwords and the number of out of vocabulary words; all these are consistently more correlated than metrics in the MeSH and CHV group. In turn, the top correlations for the HTML group, counts of P tags and list tags, were the weakest. P tags are used to create paragraphs in a Web page, being roughly a proxy for text lengthiness. Top estimators for the word frequency group are based on the Medical Reddit and PubMed counts, with correlations as high as the readability formulas. Finally, the group with the highest correlated estimators are the regressors and classifiers,



with top estimators being the Neural Network regressor and the multinomial Naive Bayes.

this section misses some short of conclusion or at least a link to the next section

## 5 Which Preprocessing Approach To Prefer

We further investigate the preprocessing steps with the groups of features introduced in Table 2. For that, we present in Figures 4 and 5 the box plot of different correlation metrics divided by preprocessing alternative for CLEF eHealth 2015 and 2016. For instance, the very first box plot in the upper part of these figures shows the absolute Pearson’s rank correlation of different readability metrics when using a combination of Naive and ForcePeriod as preprocessing steps. Boxes extend from the lower to upper quartile values of the data, with a line at the median. Whiskers extend from the box to show the range of the data. Flier points are those past the end of the whiskers, usually interpreted as outlier values.

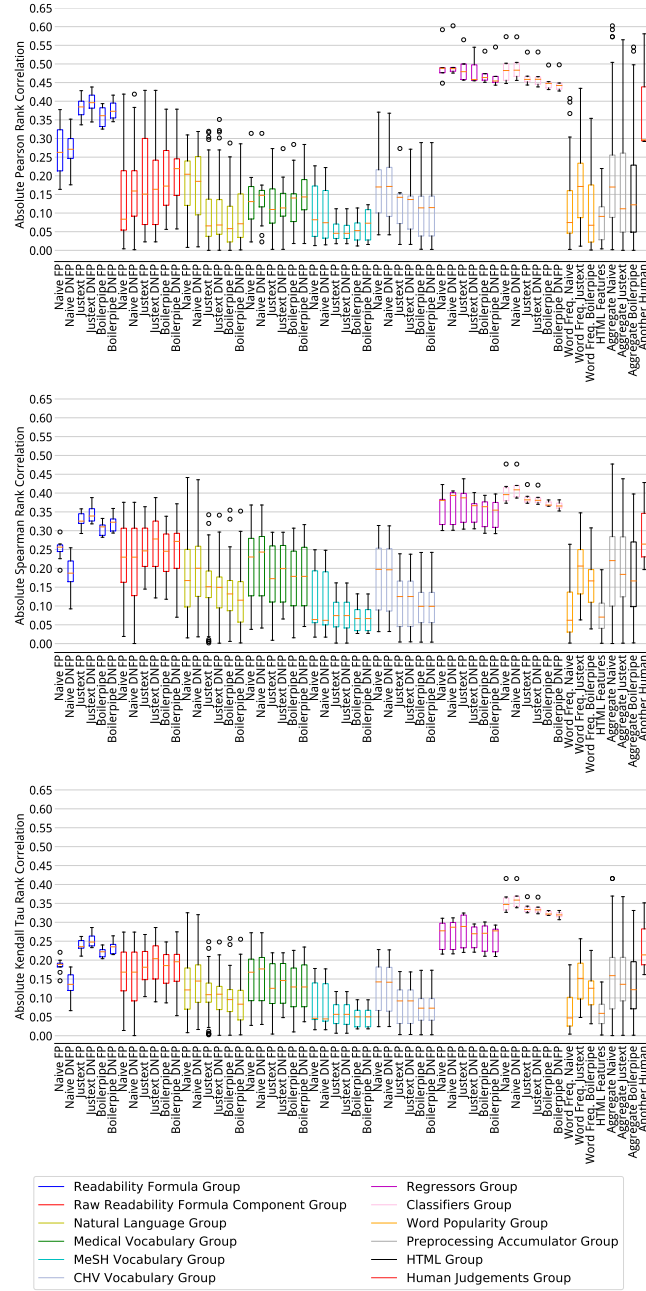
We also include in Figures 4 and 5 boxes for the summary of the 3 preprocessing procedures to remove HTML, the use of HTML features, which is done without any preprocessing and the comparison with other human assessors. For CLEF eHealth 2015, we used as human assessments the additional assessments made by unpaid medical students and health consumers (see [12]), while for CLEF eHealth 2016 data, we randomly selected 100 pages that were assessed by another assessor. **add at least another person doing assessments**. The correlations with human assessments provide important insights on how hard and subjective understandability assessments are.

Figure 4 shows the correlations for CLEF eHealth 2015 assessments. The choice of preprocessing method had the highest impact on the traditional readability formula group, with the Naive preprocessing clearly underperforming the other preprocessing methods. The choice of the Naive method was also the worst with the raw readability formula components and word frequency estimators, but, interestingly, it was a good choice, if not the best one, for all other groups. The highest correlations were archived by the regressors and classifiers, independently of the preprocessing method used.

Similarly to Figure 4, Figure 5 reports the findings for CLEF eHealth 2016. This time, though, the Naive preprocessing method was clearly underperforming for most of the groups analysed, including regressors and classifiers.

In order to further understand our experiments, we compared the median of each pair of preprocessing strategy showed in Figures 4 and 5 and present the results in Table 4. For instance, the entry *FP ; DNFP* counts the number of times the median value for ForcePeriod was superior to DoNotForcePeriod when comparisons with the same HTML cleaning method was used, e.g. Naive ForcePeriod versus Naive DoNotForcePeriod. From all comparisons, the ones that were statistically significant according to a t-test are shown inside parentheses.

The upper part of Table 4 shows results for the comparisons between ForcePeriod (FP) and DoNotForcePeriod (DNFP). Although the interpretation of readability formulas is drastically affected by this choice of preprocessing, as



**Fig. 4.** Box plots divided by feature groups. Correlations are calculated using under-standability labels from relevant documents assessed in CLEF eHealth 2015



learning in Chapter ??, the correlation results are not. The number of times FP reached a higher correlation than DNFP is roughly the same that DNFP was higher than FP. The bottom part of Table 4 shows the comparisons made for Naive, Justext and Boilerpipe. Results for CLEF 2015 contrast with 2016, while Naive was slightly better than Boilerpipe and Justext in 2015, it was the worst in almost all 2016 comparisons. Also, the comparisons between Justext and Boilerpipe are exactly the opposite from 2015 to 2016.

## 6 Experimenting with Understandability

The data use made our analysis here was also used in the Information Retrieval branch of CLEF eHealth. We focus our attention to the CLEF eHealth 2016 campaign leaving 2015 experiments offline<sup>13</sup>.

We start by the defining the evaluation measure that we will use here. In CLEF eHealth campaign, organizers used a modification of RBP which ties together the relevance of a document with any other relevance dimension, in this case in particular, with understandability. Mathematically, it consists in adding an understandability factor to the RBP formula, as shown:

add formula here

The drawback of such evaluation metric is that we cannot separately evaluate each dimension. We propose, instead, to separately evaluate a ranking list with respect to its topical relevance and its understandability:

- P@10r: a document is topically relevant if assessed as somewhat relevant or highly relevant. This metric counts the number of relevant documents in the top 10 documents of a ranked list.
- P@10u: a document is relevant for this metric if the understandability score is smaller than a threshold  $U$ . Like P@10r, we count the number of relevant documents in the first 10 docs of a ranked list. We use  $U = 40$  in our experiments.

I decided to use this threshold based on the data. I will need to add a figure to support this claim, I think.

During the campaign, organizers opt to use shallow pools and focus on highly ranked documents, using P@10r as one official metric for topical relevance. It makes our choice of metric natural. Likewise it is traditionally done with F measure, we combine P@10r and P@10u with an harmonic mean:  $F_{7u} = 2 \times \frac{P@10r \times P@10u}{P@10r + P@10u}$

## 7 Conclusion

There is an abundance of factors that affect how readability is perceived by users. In this chapter we devised and studied a large number of readability estimators, ranging from traditional readability formulas extensively used in the past 50 years to state-of-the-art machine learning algorithms. We grouped them into semantically related groups in order to visualize their correlation with human assessments collected during CLEF eHealth campaigns in 2015 and 2016.

<sup>13</sup> Link to experiments will be available upon acceptance of this manuscript

Complementary to our previous chapter, we evaluated how preprocessing steps impact the readability estimation in traditional readability formulas and in other modern estimators. We empirically learnt the importance of preprocessing steps when applying readability formulas, as the highest correlations happen when other than the Naive method is used. For the most modern estimators, such as the ones based on machine learning methods, the correlation is less sensible to the preprocessing steps.

We also studied the correlation of each individual readability formula to the human assessment to provide insights on which formula should be preferred. Our analysis concluded that the Simple Measure of Gobbledygook (SMOG) and Dale-Chall Index (DCI) were the most correlated metrics for the two datasets studied and, together with Coleman-Liau Index (CLI) and the Flesch Reading Ease (FRE) are the most stable metrics across datasets, and therefore, should be preferred.

Finally, this chapter serves as a basis for the following chapters of this work, as the learning to rank methods will largely use the estimators devised and analysed here.

## References

1. S. R. Atcherson, A. E. DeLaune, K. Hadden, R. I. Zraick, R. J. Kelly-Campbell, and C. P. Minaya. A computer-based readability analysis of consumer materials on the american speech-language-hearing association website. *Contemporary Issues in Communication Science & Disorders*, 41, 2014.
2. C. H. Björnsson. Readability of newspapers in 11 languages. *Reading Research Quarterly*, 18(4):480–497, 1983.
3. M. Coleman and T. L. Liau. A Computer Readability Formula Designed for Machine Scoring. *Journal of Applied Psychology*, 1975.
4. K. Collins-Thompson. Computational assessment of text readability: A survey of current and future research. *ITL-International Journal of Applied Linguistics*, 165(2):97–135, 2014.
5. E. Dale and J. S. Chall. A formula for predicting readability: Instructions. *Educational Research Bulletin*, 27(2):37–54, 1948.
6. P. Fitzsimmons, B. Michael, J. Hulley, and G. Scott. A readability assessment of online parkinson’s disease information. *The journal of the Royal College of Physicians of Edinburgh*, 40(4):292–296, 2010.
7. M. A. Graber, C. M. Roller, and B. Kaeble. Readability levels of patient education material on the world wide web. *Journal of Family Practice*, 48(1):58–59, 1999.
8. R. Gunning. *The Technique of Clear Writing*. McGraw-Hill, 1952.
9. J. Kincaid, R. Fishburne, R. Rogers, and B. Chissom. *Derivation of New Readability Formulas for Navy Enlisted Personnel*. National Technical Information Service, 1975.
10. G. H. McLaughlin. SMOG Grading - a New Readability Formula. *Journal of Reading*, 1969.
11. A. Meillier and S. Patel. Readability of healthcare literature for gastroparesis and evaluation of medical terminology in reading difficulty. *Gastroenterology Research*, 10(1):1–5, 2017.

12. J. Palotti, G. Zuccon, J. Bernhardt, A. Hanbury, and L. Goeuriot. Assessors Agreement: A Case Study across Assessor Type, Payment Levels, Query Variations and Relevance Dimensions. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 7th International Conference of the CLEF Association, CLEF'16 Proceedings*. Springer International Publishing, 2016.
13. J. Palotti, G. Zuccon, L. Goeuriot, L. Kelly, A. Hanbury, G. J. F. Jones, M. Lupu, and P. Pecina. ShARe/CLEF eHealth Evaluation Lab 2015, Task 2: User-centred Health Information Retrieval. In *Working Notes for CLEF 2015 Conference, Toulouse, France, September 8-11, 2015.*, 2015.
14. C. R. Patel, D. V. Cherla, S. Sanghvi, S. Baredes, and J. A. Eloy. Readability assessment of online thyroid surgery patient education materials. *Head & neck*, 35(10):1421–1425, 2013.
15. E. A. Smith and R. J. Senter. *Automated Readability Index*. AMRL-TR-66-220. Aerospace Medical Research Laboratories, 1967.
16. L. Soldaini, A. Cohan, A. Yates, N. Goharian, and O. Frieder. *Retrieving Medical Literature for Clinical Decision Support*, pages 538–549. Springer International Publishing, 2015.
17. R. C. Wiener and R. Wiener-Pla. Literacy, pregnancy and potential oral health changes: The internet and readability levels. *Maternal and child health journal*, 18(3):657–662, 2014.
18. X. Yan, R. Y. Lau, D. Song, X. Li, and J. Ma. Toward a semantic granularity model for domain-specific information retrieval. *ACM Transactions on Information Systems*, 29(3):15:1–15:46, July 2011.
19. Q. T. Zeng and T. Tse. Exploring and developing consumer health vocabularies. *Journal of the American Medical Informatics Association*, 13(1):24–29, 2006.
20. W. Zhou, V. Torvik, and N. Smalheiser. Adam: Another database of abbreviations in medline. *Bioinformatics*, 22(22):2813–2818, 2006.
21. G. Zuccon, J. Palotti, L. Goeuriot, L. Kelly, M. Lupu, P. Pecina, H. Mueller, J. Budaheer, and A. Deacon. The IR Task at the CLEF eHealth evaluation lab 2016: user-centred health information retrieval. In *CLEF 2016-Conference and Labs of the Evaluation Forum*, volume 1609, pages 15–27, 2016.

**Table 2.** Metrics used as understandability proxies;  $\star$ : raw values are used.  $\diamond$ : values normalised by number of words in a documents are used.  $\dagger$ : values normalised by number of sentences in a document are used.

Group	Metric	Group	Metric
Traditional Readability Formulas	Automated Readability Index (ARI) [15]	HTML Features	# of Abbr tags
	Coleman-Liau Index (CLI) [3]		# of A tags
	Dale Chall Index (DCI) [5]		# of Blockquote tags
	Flesch-Kincaid Grade Level (FKGL) [9]		# of Bold tags
	Flesch Reading Ease (FRE) [9]		# of Cite tags
	Gunning Fog Index (GFI) [8]		# of Div tags
	Lasarhetsindex (LIX) [2]		# of Form tags
Raw Components of Readability Measures	Simple Measure of Gobbledygook (SMOG) [10]		# of H1 tags
	# of Characters $\star\diamond\dagger$		# of H2 tags
	# of Words $\star\diamond$		# of H3 tags
	# of Sentences $\star\diamond$		# of H4 tags
	# of Difficult Words (Dale Chall list [5]) $\star\diamond\dagger$		# of H5 tags
	# of Words Longer than 4 chars $\star\diamond\dagger$		# of H6 tags
	# of Words Longer than 6 chars $\star\diamond\dagger$		# of Hs (any H above)
	# of Words Longer than 10 chars $\star\diamond\dagger$		# of Img tags
	# of Words Longer than 13 chars $\star\diamond\dagger$		# of Input tags
	# of Number of Syllables $\star\diamond\dagger$		# of Link tags
Medical Vocabularies	# of Polysyllable Words ( $\geq 3$ Syllables) $\star\diamond\dagger$	Word Frequency	# of DL tags
	# of Words with Medical Prefix $\star\diamond\dagger$		# of UL tags
	# of Words with Medical Suffix $\star\diamond\dagger$		# of OL tags
	# of Acronyms $\star\diamond\dagger$		# of List (DL + UL + OL)
	# of ICD Concepts $\star\diamond\dagger$		# of Q tags
	# of Drugbank $\star\diamond\dagger$		# of Scripts tags
Consumer Health Vocabulary (CHV) [19] Features	# of Words in medical dict. (OpenMedSpel) $\star\diamond\dagger$		# of Spans tags
	CHV Mean Score for all Concepts $\star\diamond\dagger$		# of Table tags
	# of CHV Concepts $\star\diamond\dagger$		# of P tags
	CHV Mean Score for Symptom Concepts $\star\diamond\dagger$	Regressor	25th percentil English Wikipedia
	# of CHV Symptom Concepts $\star\diamond\dagger$		50th percentil English Wikipedia
Medical Subject Headers (MeSH)	CHV Mean Score for Disease Concepts $\star\diamond\dagger$		75th percentil English Wikipedia
	# of CHV Disease Concepts $\star\diamond\dagger$		Mean Rank English Wikipedia
	# of MeSH Concepts $\star\diamond\dagger$		Mean Rank English Wikipedia - Includes OV
	Average Tree of MeSH Concepts $\star\diamond\dagger$		25th percentil Medical Reddit
	# of MeSH Symptom Concepts $\star\diamond\dagger$		50th percentil Medical Reddit
Natural Language	Average Tree of MeSH Symptom Concepts $\star\diamond\dagger$		75th percentil Medical Reddit
	# of MeSH Disease Concepts $\star\diamond\dagger$		Mean Rank Medical Reddit
	Average Tree of MeSH Disease Concepts $\star\diamond\dagger$		Mean Rank Medical Reddit include OV
	Positive Words $\star\diamond\dagger$		25th percentil Pubmed
	Negative Words $\star\diamond\dagger$		50th percentil Pubmed
	Neutral Words $\star\diamond\dagger$		75th percentil Pubmed
	# of verbs $\star\diamond\dagger$		Mean Rank Pubmed
	# of nouns $\star\diamond\dagger$		Mean Rank Pubmed - Includes OV
	# of pronouns $\star\diamond\dagger$		25th p. Wikipedia+Reddit+Pubmed
	# of adjectives $\star\diamond\dagger$		50th p. Wikipedia+Reddit+Pubmed
	# of adverbs $\star\diamond\dagger$	Classifier	75th p. Wikipedia+Reddit+Pubmed
	# of adpositions $\star\diamond\dagger$		Mean R. Wiki.+Reddit+Pubmed
	# of conjunctions $\star\diamond\dagger$		Mean R. Wiki.+Reddit+Pubmed - w. OV
	# of determiners $\star\diamond\dagger$		Linear Regressor
	# of cardinal numbers $\star\diamond\dagger$		Gradient Boosting Regressor
	# of particles or other function words $\star\diamond\dagger$		Multi-layer Perceptron Regressor
	# of other POS (foreign words, typos) $\star\diamond\dagger$		Random Forest Regressor
	# of punctuation $\star\diamond\dagger$		Support Vector Machine Regressor
	Height of part-of-speech parser tree $\star\diamond\dagger$		Logistic Regression
	# of Entities $\star\diamond\dagger$		Gradient Boosting Classifier
	# of Stopwords $\star\diamond\dagger$		Multinomial Naive Bayes
	# of words not found in Aspell Eng. dict. $\star\diamond\dagger$		Multi-layer Perceptron Classifier
			Random Forest Classifier
			Support Vector Machine Classifier

**Table 3.** Metrics with highest correlation per group. In bold are the metric that archived the highest correlation for a correlation measure.

Dataset	Group	Metric	Preprocessing	Pearson	Spearman	KendallTau
CLEF 2015	Readability Formula	SMOG Index	Justext NFP	<b>0.438</b>	<b>0.388</b>	<b>0.286</b>
	Formula Component	Avg. Number of Polysyl. Words per Word	Justext FP	<b>0.429</b>	0.364	0.268
		Avg. N. of Polysyl. Words per Sentence	Justext NFP	0.192	<b>0.388</b>	<b>0.286</b>
	Medical Vocabulary	Avg. N. Medical Prefixes per Word	Naive FP	<b>0.314</b>	0.312	0.229
		Number of Medical Prefixes	Naive FP	0.131	<b>0.368</b>	<b>0.272</b>
	CHV	CHV Mean Score for all Concepts	Naive FP	<b>0.371</b>	<b>0.314</b>	<b>0.228</b>
	MeSH	Number of MeSH Concepts	Naive FP	<b>0.227</b>	<b>0.249</b>	<b>0.178</b>
	Natural Language	N. of words not found in Aspell Dict. Number of Pronouns per Word	Justext NFP	<b>0.351</b>	0.276	0.203
			Naive FP	0.271	<b>0.441</b>	<b>0.325</b>
	HTML	Number of P Tags	None	<b>0.219</b>	<b>0.196</b>	<b>0.142</b>
CLEF 2016	Readability Formula	Dale Chall Index Dale Chall Index	Justext FP	<b>0.439</b>	0.381	0.264
			Boilerp. FP	0.437	<b>0.382</b>	<b>0.264</b>
	Formula Component	Avg. Difficult Words Per Word	Boilerp. FP	<b>0.431</b>	<b>0.379</b>	<b>0.262</b>
	Medical Vocabulary	Avg. Prefixes per Sentence ICD Concepts Per Sentence	Justext FP	<b>0.263</b>	0.242	0.164
			Justext NFP	0.014	<b>0.253</b>	<b>0.172</b>
	CHV	CHV Mean Score for all Concepts CHV Mean Score for all Concepts	Justext FP	<b>0.329</b>	0.313	0.216
			Boilerp. FP	0.329	<b>0.325</b>	<b>0.224</b>
	MeSH	Number of MeSH Concepts Number of MeSH Disease Concepts	Boilerp. NFP	<b>0.201</b>	0.166	0.113
			Boilerp. NFP	0.179	<b>0.192</b>	<b>0.132</b>
	Natural Language	Avg. Stopword Per Word Number of Pronouns	Boilerp. FP	<b>0.344</b>	0.312	0.213
			Boilerp. FP	0.341	<b>0.364</b>	<b>0.252</b>
CLEF 2016	HTML	Number of Lists Number of P Tags	None	<b>0.114</b>	0.021	0.015
				0.110	<b>0.123</b>	<b>0.084</b>
	Word Frequency	Mean Rank Medical Reddit 50th percentil Medical Reddit	Boilerp. NFP	<b>0.387</b>	0.312	0.214
			Justext NFP	0.351	<b>0.315</b>	<b>0.216</b>
	Regressors	Neural Network Regressor Random Forest Regressor	Justext NFP	<b>0.454</b>	<b>0.373</b>	0.258
			Boilerp. NFP	0.389	0.355	<b>0.264</b>
	Classifiers	Multinomial Naive Bayes	Justext FP	<b>0.461</b>	<b>0.391</b>	<b>0.318</b>

**Table 4.** Exhaustive Comparison summary using the data from Figures 1.2 and 1.3. Numbers inside parentheses represent the number of tests that yielded  $p \leq 0.05$  in a two-tailed t-test

Comparison	CLEF 2015				CLEF 2016			
	Pearson	Spearman	Kendall	Tau Total	Pearson	Spearman	Kendall	Tau Total
FP $\neq$ DNFP	8 (0)	11 (4)	11 (3)	30 (7)	16 (10)	10 (3)	11 (4)	37 (17)
FP $\leq$ DNFP	16 (5)	12 (5)	12 (6)	40 (16)	8 (0)	12 (2)	11 (2)	31 (4)
FP == DNFP	0	1	1	2	0	2	2	4
Naive $\neq$ Justext	11 (7)	9 (6)	9 (5)	29 (18)	1 (0)	0 (0)	0 (0)	1 (0)
Naive $\leq$ Justext	6 (4)	8 (4)	8 (4)	22 (12)	16 (12)	17 (13)	17 (13)	50 (38)
Naive == Justext	0	0	0	0	0	0	0	0
Naive $\neq$ Boilerpipe	12 (7)	10 (6)	10 (5)	32 (18)	0 (0)	0 (0)	0 (0)	0 (0)
Naive $\leq$ Boilerpipe	5 (4)	7 (3)	7 (3)	19 (10)	16 (12)	17 (13)	17 (13)	51 (39)
Naive == Boilerpipe	0	0	0	0	0	0	0	0
Justext $\neq$ Boilerpipe	10 (7)	16 (9)	14 (8)	40 (24)	9 (4)	9 (4)	4 (2)	17 (8)
Justext $\leq$ Boilerpipe	7 (2)	1 (0)	3 (1)	11 (3)	8 (2)	8 (2)	13 (2)	34 (6)
Boilerpipe == Justext	0	0	0	0	0	0	0	0



**Table 5.** Reports on the experiments with 4 base runs: the top 3 runs of CLEF eHealth 2017 and a plain baseline run. The second (indices 5-8) and third (indices 9-12) parts shows results of a RegressorTree based on top features from Table X. The forth part of this table shows results when reranking top 20 results based on Dale-Chall Index. Finally, we combine selected runs with reciprocal rank fusion in the last part of this table. Results shows that regression on top 15 improves understandability to the detriment of topical relevance. Regression on top 20 it is even more aggressive. Dale Chall presents the same logic, but understandability gains are smaller compared to regressor results. Finally, the combination with rrf yields the best  $F_{ru}$  scores. I still need to compute if we get statistically significant improvements or not.

Index	Rerank	Run	$P_r@10$	$P_u@10$	$F_{ru}$	RBP	uRBP	Unj@10	$P_r@10$	$P_u@10$	$F_{ru}^*$
1	No Rerank	BM25 Q.E.	31.43	49.70	38.51	29.05	18.65	0.02	31.70	50.83	39.05
2		GUIR	29.67	50.97	37.50	28.11	18.12	0.01	30.20	51.80	38.16
3		ECNU	29.33	49.27	36.77	27.70	17.55	0.01	29.50	49.87	37.07
4		Plain BM25	26.47	46.33	33.69	25.28	16.05	0.06	27.63	48.93	35.32
5	Regress Top 15	Based on Run 1	27.87	49.40	35.63	26.38	19.77	0.01	28.10	49.60	35.88
6		Based on Run 2	27.27	48.57	34.93	25.95	19.31	0.01	27.30	48.60	34.96
7		Based on Run 3	26.60	50.10	34.75	25.08	18.68	0.00	26.67	50.10	34.81
8		Based on Run 4	23.73	50.87	32.37	22.86	17.68	0.09	24.97	52.67	33.87
9	Regress. Top 20	Based on Run 1	27.70	<b>52.13</b>	36.18	26.10	20.06	0.02	28.17	52.53	36.67
10		Based on Run 2	26.70	<b>53.10</b>	35.53	25.44	19.62	0.01	26.83	53.33	35.70
11		Based on Run 3	26.17	<b>52.67</b>	34.96	24.88	19.21	0.02	26.43	53.17	35.31
12		Based on Run 4	23.30	<b>51.43</b>	32.07	22.53	18.01	0.12	25.10	55.20	34.51
13	Dale-Chall top 20	Based on Run 1	29.63	49.93	37.19	26.45	17.78	0.02	29.87	50.37	37.50
14		Based on Run 2	28.37	51.10	36.48	25.01	16.96	0.01	28.60	51.40	36.75
15		Based on Run 3	27.70	50.80	35.85	24.54	16.49	0.02	28.03	51.47	36.30
16		Based on Run 4	24.13	49.97	32.55	21.28	14.60	0.10	26.17	53.17	35.07
17	Reciprocal Rank Fusion	Runs 1 and 9	30.70	52.37	<b>38.71</b>	27.93	20.20	0.02	31.10	53.30	39.28
18		Runs 2 and 10	30.67	51.90	<b>38.55</b>	28.34	20.11	0.01	30.97	52.53	38.96
19		Runs 3 and 11	30.40	51.37	<b>38.20</b>	27.82	19.72	0.01	30.67	51.90	38.55
20		Runs 4 and 12	28.67	50.67	<b>36.62</b>	26.33	18.68	0.03	29.37	51.87	37.50