

# SW4ALL: a CEFR-Classified and Aligned Corpus for Language Learning

Rodrigo Wilkens, Leonardo Zilio and Cédric Fairon

Centre de Traitement Automatique du Langage - CENTAL

Université catholique de Louvain - UCL

{rodrigo.wilkens, leonardo.zilio, cedrick.fairon}@uclouvain.be

## Abstract

Learning a second language is a task that requires a good amount of time and dedication. Part of the process involves the reading and writing of texts in the target language, and so, to facilitate this process, especially in terms of reading, teachers tend to search for texts that are associated to the interests and capabilities of the learners. But the search for this kind of text is also a time-consuming task. By focusing on this need for texts that are suited for different language learners, we present in this study the SW4ALL, a corpus with documents classified by language proficiency level (based on the CEFR recommendations) that allows the learner to observe ways of describing the same topic or content by using strategies from different proficiency levels. This corpus uses the alignments between the English Wikipedia and the Simple English Wikipedia for ensuring the use of similar content or topic in pairs of text, and an annotation of language levels for ensuring the difference of language proficiency level between them. Considering the size of the corpus, we used an automatic approach for the annotation, followed by an analysis to sort out annotation errors. SW4ALL contains 8,669 pairs of documents that present different levels of language proficiency.

**Keywords:** SLA, CEFR classification, Aligned corpus

## 1. Introduction

Learning a second language is a process that requires exposition to texts, especially for the acquisition of vocabulary (Rott, 1999). To retrieve texts that match learners' language level (or proficiency) it is possible to use a corpus carefully designed for language learners, or one can search for them in the web. Following these alternatives, systems can dig up texts aiming at finding those that are best suited to the language skills of a learner. Examples of these systems are REAP (Heilman et al., 2008), FLAIR (Chinkina et al., 2016), and READ-X (Miltakaki and Troutt, 2007), which use the web as a corpus. The use of web allow the learners to interact with a huge diversity of texts, which makes it easier to find those that correspond to their interests. But most of the web texts retrieved by search engines require a high language proficiency, even for native speakers (Vajjala and Meurers, 2013). On the other hand, the use of an off-line corpus ensures content quality, while hindering the search for different text topics.

This dichotomy of text sources enforces different types of restrictions on the systems. An alternative for trying to get the best of both approaches is to use the web as source of texts, but restricting it to trustful domains. This is similar to SourceFinder's method, in which on-line newspapers and magazines are downloaded and processed as text sources (Passonneau et al., 2002; Sheehan et al., 2007). However, this kind of approach doesn't allow for an easy update of the texts, because the content is stored off-line, and an update would require a rerun of the whole corpus compilation process. Another reliable type of text source that is adequate to language learners is corpora made up of simplified texts, such as the Weekly Reader, the Simple English Wikipedia, and the BBC Bitesize. This type of source generally represents an attempt to make texts more accessible, by adapting or simplifying them to present a language that should be easier to understand for a non-native speaker. Wikipedia also has concerns on the comprehension abili-

ties of its users, so that, for the English language, there is a simplified version, addressing the needs of natives with low literacy, but also the needs of learners of English. This type of resource presents a simplified version of a source article, serving as a facilitator for the communication of knowledge for those with less language skill, but it doesn't present an information about to what extent the text is simplified. It doesn't explain the simplification strategies applied to each text or for what target reader each text was simplified, and this, in the case of language learning, which categorizes learners in different levels, is crucial to better inform the learner about which texts would be at an understandable level. For instance, the Common European Framework of Reference for Languages (CEFR) classifies learners in six language levels (ranging from A1 to C2), while other frameworks, like the Cambridge ESOL classifies them in 5 levels, and still other frameworks use score ranges, like TOEFL and IELTS. Without this information about the language level, or some other information about the adequacy to a given target reader, language resources that present a simplified version are, per se, not well suited for language learners, because the simplification requirements may not be of use for the needs of specific learners. As such, another layer with a more pedagogically relevant classification is needed.

This study aims at automatically determining the CEFR level of pairs of original and simplified texts, so that a corpus of paired texts pertaining to different language levels can be used in a language learning framework.<sup>1</sup> This would allow learners to compare text structures from different levels that describe the same content, while also allowing for the selection of topics of interest. To this end, we annotated language levels in an aligned version of the English

---

<sup>1</sup>The developed resource is available at [http://cental.uclouvain.be/resources/smalla\\_smille/sw4all/](http://cental.uclouvain.be/resources/smalla_smille/sw4all/)

Wikipedia<sup>2</sup> (EW) and Simple English Wikipedia<sup>3</sup> (SEW), and filtered pairs of texts that are associated to different levels, but that refer to the same topic or content. This process resulted in a resource that we called *Simple Wikipedia for Aligned Language Learning* (SW4ALL).

This paper is organized as follows: Section 2. presents systems that classify texts for the purpose of language learning; Section 3. describes the training corpus, the detailed training methodology and its results, and the application of the trained model to the aligned EW-SEW; Section 4. is where the evaluation of the classification model is presented; Section 5. contains a description of the resulting annotated corpus; and, finally, Section 6. is reserved for our final remarks on this study.

## 2. Related Work

The search of texts that can improve language learning skills, and, at the same time, be able to match the learners' interests is a very time-consuming activity for language teachers. Aiming to reduce the time allocated to this task, the SourceFinder (Passonneau et al., 2002; Sheehan et al., 2007) allows teachers to search for documents classified in different language levels (according to the Graduate Record Examinations curriculum) by means of keywords. One of the advantages of SourceFinder is the use of off-line texts, which allows the processing of the texts with several NLP tools without delay to the user. On the other hand, it only allows the search for text content and grammar structures. Using online documents, the REAP (Heilman et al., 2008), READ-X (Miltakaki and Troutt, 2007), and LAWSE (Ott and Meurers, 2011) systems allow the users (learners or teachers) to search texts by means of keywords and to filter them according to readability measures. Those systems identify the text readability by applying traditional readability scores (Flesch-Kincaid measure (Kincaid et al., 1975), and Gunning Fog Index (Gunning, 1952)) that are based on shallow cues (e.g. number of words per sentences and syllables per word). These measures have the advantage of a fast annotation process, but they are not accurate, and they require the users to deal with a score that may not be familiar to them.

Using a more accurate text classification method, the FLAIR system (Chinkina et al., 2016; Chinkina and Meurers, 2016) dynamically crawls, annotates, and classifies the 20 first results of a search engine. The FLAIR text classification is based on parsing information and on the official English language curriculum of schools in Baden-Württemberg (Chinkina et al., 2016).

Taking into account the pedagogical function of these systems, a major point is their ability to address documents that are readable by learners. However, text length-based readability scores weigh only sentence length and word difficulty, ignoring factors such as cohesion (Bruce et al., 1981). Recently, Xia et al. (2016) compared syntactic and length-based features for text classification according to language level, and identified that adding syntactic features on top of length-based features improved the classification results,

but using only length-based features presented a better result than the syntactic features alone.

Regarding the representation of the texts in features, there is a huge variety of options in the literature. They can be grouped into 6 categories: length-based (e.g. word and sentence length), lexical (e.g. proportion of words in a list of easy words), morphological (e.g. part of speech), syntactic (e.g. presence of passive voice), semantic (e.g. word polysemy), and language model (e.g. n-gram model perplexity). The syntactic features could be split into two groups, depending on how the parser is used. Usually, a parser-based annotation of features follows the same process as the morphological annotation: simple counts of parser annotation. However, some studies, such as François and Fairon (2012) and Callan and Eskenazi (2007), used information beyond parsed tags.

All those systems focus on presenting a readable text, but some of them go beyond that, presenting exercises for supporting the learning activity in a more active way (e.g. REAP). In spite of all the effort to present readable and interesting texts to learners, those systems do not indicate how learners can improve their skills by using the indicated texts.

## 3. Methodology

Regarding the objective of automatically determining pairs of texts that present good examples of different language levels, we trained a classifier and applied it to the aligned English Wikipedia (EW) and Simple English Wikipedia (SEW). This annotated resource was named *Simple Wikipedia for Aligned Language Learning* (SW4ALL). In this section, we present first the alignment between the two Wikipedia versions (Section 3.1.). We then move on to the resources needed to build the classifier: the training corpus (Section 3.2.) and the feature set (Section 3.3.). Finally, we discuss the application of the classifier to the aligned EW-SEW (Section 3.4.).

### 3.1. Aligned Wikipedias

The Wikipedia is a collaborative encyclopedia with huge amounts of texts available in several languages. In English, there are two version, one that focus just on encyclopedic information, and the other that requires the content to be written in a simplified way. Comparing the vocabulary of the two encyclopedias, Coster and Kauchak (2011) identified that 96% of the words in the simple version are found in the other version, and 87% of the words in the normal version are found in the simple version. This overlap is also found at the n-gram level. Regarding the alignment of the Wikipedias, there are different versions, and in this paper we opted for the version organized by Kauchak (2013), in which the texts were aligned both at the document and sentence level. Kauchak (2013) downloaded and cleaned all articles from the Wikipedias (removing stubs and navigation pages), resulting in 60K articles each. The difference in the number of sentences between the Wikipedia versions is partially because some articles from SEW present just partial information. Indeed, the sentence level alignment, also presented by Kauchak (2013), was possible in only 28% of

---

<sup>2</sup><https://wikipedia.org/>

<sup>3</sup><https://simple.wikipedia.org/>

the sentences from the SEW (and in 4.25% of the sentences from EW).

### 3.2. Training Corpus

Focusing on improving the learners' writing skills, we opted to use a corpus of texts written by language learners. In that way we are able to present texts compatible with learners' productive skills, making it easier for them to identify the structures that exist in those texts. So, by applying a model that was trained on texts produced by learners to the Wikipedias, we expect that the structures in the text will be familiar to the language learners, while also providing an authentic source of information, because the texts of the Wikipedias were written by native speakers.

In this study, we used the EF-Cambridge Open Language Database (EFCAMDAT) (Geertzen et al., 2013) as training corpus. This corpus is divided according to the Common European Framework of Reference for languages (CEFR) (Verhelst et al., 2009), containing a total of 532 thousand documents (33 million tokens) written by 83,385 learners of 137 nationalities, and each document has an evaluation score and an associated topic (e.g. introducing yourself by email). The data is distributed into three main levels, each with two sublevels (all referenced by a letter and a number): basic (breakthrough or A1; waystage or A2), independent (threshold or B1; vantage or B2), and proficient (effective operational proficiency or C1; mastery or C2). The EFCAMDAT corpus contains an unbalanced number of documents per level (e.g. 151 thousand documents for A2 and 23 thousand for B2), so we selected 9,000 random documents from each level, and also filtered out those documents that did not achieve an evaluation score higher than 80%, because, in these cases, learners' errors could have an impact on the machine learning approach (Pilán et al., 2016). The result of this process is a corpus of 40,946 documents (9,000 for levels A1, A2, B1, and B2, and 4,946 for C1<sup>4</sup>).

### 3.3. Feature Annotation

The annotation process was three-fold: first the documents were automatically parsed with the Stanford Parser (Manning et al., 2014), and then a series of features were annotated, including the ones developed by project SMILLE (Zilio and Fairon, 2017), which have a good performance, comparable to state-of-the-art parsers in the labeling task (Zilio et al., 2017a; Zilio et al., 2017b). Finally, we annotated the documents with readability scores. The annotations were grouped in four categories, inspired by Xia et al. (2016): length-based, morphological, syntactic, and readability.

In addition to these sets of features, we also took into account the grouping of these morphological and syntactic features according to two criteria: CEFR level (e.g., A1, A2, etc.) in which they should be learned<sup>5</sup>, which resulted

<sup>4</sup>We used all documents that were scored over 80% C1 data, and we did not use the C2-level documents due to the low number of documents.

<sup>5</sup>For allocating each structure to a given CEFR level, we used SMILLE's pedagogical model.

in 5 grammatical and 5 word features; and type of grammatical structure (still respecting the CEFR levels division; for instance, connectives, which are learned on CEFR level B1, were put together in one set, but modals, which are learned in different levels, were separated in two sets). These sets of features were called pedagogical feature sets.

### 3.4. Annotating the Aligned EW-SEW

The model trained on the Cambridge Corpus data was applied to the aligned version of the English Wikipedia (EW) and the Simple English Wikipedia (SEW), so that we could observe which pairs of texts are suitable for contrasting different CEFR levels of English. Based on the premises of the SEW that the texts should be simpler than the EW, they should at least be on the same CEFR level as their EW counterpart, so we considered that all pairs for which the system classified the SEW text as having a higher CEFR level than the EW were bad for SW4ALL. Conversely, all the pairs for which the system presented the SEW text as being from a lower CEFR level than its EW counterpart were considered good for the resource. Pairs which presented the same level for both Wikipedias were further analyzed for level tendencies, as we discuss in Section 5..

## 4. Model evaluation

The first topic to be addressed in this evaluation is the quality of the model used to annotate the corpus. First, in Section 4.1., we evaluate the prediction power of each feature. Then, we discuss the corpus size impact on the model's performance (Section 4.2.).

### 4.1. Feature selection

To compare the features' quality, we ranked all of them according to the Gain Ratio algorithm (Frank et al., 2009), an entropy-based feature selection algorithm which ranks each feature according to its pertinence for separating the classes. We observed that 5 pairs of features (2 morphological, 2 syntactic and 6 pedagogical) presented a low score difference ( $< 0.00001$ ) in the pair. This happens because the value of each pair comes mostly from one feature, so the group effect is not observed. As such, we changed the rank position of these features to bear the same value of those with which there was no substantial difference, aiming to make a fairer comparison. We also removed 12 features scored as zero by the Gain Ratio algorithm<sup>6</sup>.

The rank distribution of the different types of features is presented in Figure 1, which displays the overall importance of pedagogical features for the model. However, it is important to notice that the rank is topped by length-based features, followed by morphological features. Interestingly, in the top 10 features, 5 are pedagogical, and the best of them is grammar-based, while the others are all vocabulary-based. Since the rank distribution does not seem to follow a normal distribution, we looked at median and quartiles of the feature types. Ranking them, we observed this, in ascending order: pedagogical, morphological, length-based,

<sup>6</sup>From the features scored as zero, one is length-based, one is vocabulary-based, four are morphological, four are syntactic, and two are pedagogical.

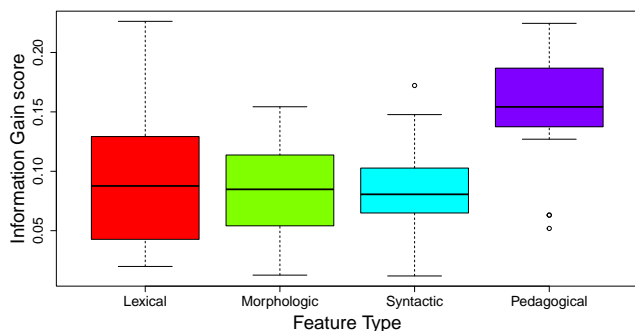


Figure 1: Score of each type of feature according to the Information Gain algorithm

syntactic, and readability-based features. Considering the rank distribution of the features, we observed that the Q1 (top 25% of the features) was ranked lower than 16.25 for length-based, 23.25 for pedagogical, 23 for morphological, 47.75 for syntactic, and 66 for readability features.

The results suggest that length-based features have great relevance for the level classification, followed by the pedagogical features. They also indicate no large difference between the morphological and syntactic features. The result of the pedagogical features is not a surprise, since they are the association of morphological and syntactic features to a pedagogical curriculum.

#### 4.2. Corpus size

Along with the feature weight, the corpus size is an important feature in machine learning. So, to evaluate the quality increase of our model in relation to corpus size, we performed ten experiments with varying corpus sizes, from 10% to 100% of the corpus. In each test, we used all feature sets and performed a ten-fold cross-validation with the Simple Logistic and the Random Forest algorithms, as can be seen in Figure 2. The model performance was evaluated in terms of precision, recall and f-measure.

The f-measure increases an average of 0.15% for each increment of 10% in corpus size. However, in regard to statistical confidence, we identify a significant increase of f-measure only when the corpus is increased by at least 20% (1,590 instances), and no difference was observed in sizes larger than 40%. Despite the nonexistence of statistical difference in larger samples, they present a smaller standard deviation. In other words, the result is more reliable using larger samples.<sup>7</sup>

To the best of our knowledge, there are no studies addressing the classification of texts written by English learners. However, in the literature there are some studies that are similar to ours. For instance, Pilán et al. (2016) address the same task, but using the MERLIN corpus (Wisniewski et al., 2013), which contains documents written in Czech, German, and Italian (80% of f-measure), and Xia et al. (2016) employ a similar set of features, but using the Cambridge English Exams dataset, which is made up of text

<sup>7</sup>Analyzing the results of the model, even with a larger corpus, we expect that a similar performance should be achieved.

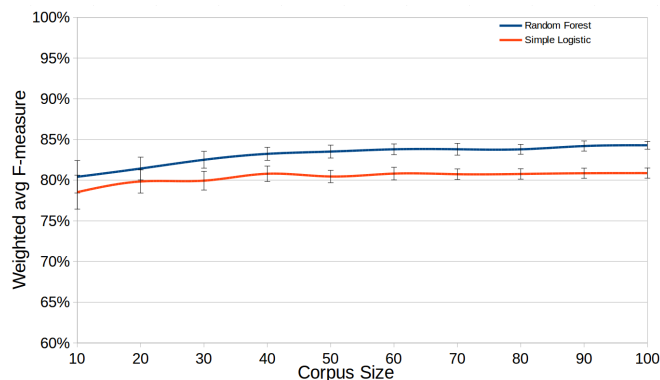


Figure 2: Average f-measure and standard deviation for the ten-fold cross-validation of the Simple Logistics and Random Forest models

written by native speakers (80% of accuracy). Still, in an effort to establish a basis of comparison, using corpus sizes that were similar to those studies, we achieved an F1 of 82% and an accuracy of 80%. If we consider the full corpus, we have an F1 of 84.7%.

### 5. Results

For developing a CEFR-classified corpus, we annotated the aligned texts of the English Wikipedia (EW) and the Simple English Wikipedia (SEW). Applying a conservative approach, we considered it a classification error when the pairs of documents had a lower language level annotated for the EW. From more than 60 thousand pairs of texts, a good amount (10,225) was classified as having the same level in both Wikipedias. For these, we further looked at the distribution as a tie breaker. For instance, a pair in which both texts were classified as B1 was investigated to see if the distribution tended to show that the SEW text was easier than the text from the EW. This process identified 2,223 pairs of documents for which the SEW version tends to present a lower level. The system classified 9,222 pairs as having an SEW document that was classified as at least one level lower than its EW counterpart, which forms a more reliable set of pairs for language learning purposes. This process left us with 11,445 pairs of texts in which the SEW document was deemed to present a lower level in relation to its EW counterpart.

To ensure the quality of the resource, we turned ourselves once again to the SEW assumptions, which indicate that texts should explain complex concepts to the user, while also splitting complex sentences, so as to form shorter, simpler sentences for the reader, but presenting the same content and with even longer texts (due to the explanations). With this information in mind, we cleaned from the corpus pairs in which the SEW text had a size (in number of words) of 90% or less than its EW counterpart, for the pair would almost certainly not present the same content, let alone a SEW version with more explicitations<sup>8</sup>. This cleaning process removed 1,359 pairs of documents from

<sup>8</sup>We did not restrict a maximum size, because it would be impractical to establish how much explicitations or sentence splitting would be too much.

the good sample, resulting in a subcorpus from the aligned EW-SEW containing 10,086 documents.

As a final step for ensuring the reliability of our classification for a user of the resource, we clustered the distribution of probabilities for each level from the classifier using the k-means<sup>9</sup> algorithm Arthur and Vassilvitskii (2007), so as to distinguish how well our classified data matched the assumptions of the SEW. We organized the clusters in suited or non-suited according to the purity as a measure of confidence. We were able to identify documents that probably present labeling error from the annotator and documents that are correctly classified. Considering only those documents that had a confidence score of more than 95%, SW4ALL consists of 6,394 pairs of documents (63% of all the documents that were considered suited), but, by relaxing this confidence to all of those above 85%, the size of the resource increases to 8,669 pairs of documents (86% of the documents that were considered suited), while maintaining a good confidence in the classification.

## 6. Conclusion

In this paper, we presented SW4ALL, a resource that focuses on the contrast of aligned texts that belong to different CEFR levels. This resource could be employed by teachers or students to compare grammar, vocabulary and general text structure of texts in different levels, but with roughly the same content.

A classification model was trained using an annotated version of the EFCAMDAT corpus, and the model was then applied to classify pairs of aligned documents from the English Wikipedia (EW) and its simplified version, the Simple English Wikipedia (SEW). After further analysis, pairs of documents in which the document from the SEW were classified as having a lower level, or a tendency to have a lower level, were used in SW4ALL, resulting in a total of 8,669 pairs of documents.

The pairs of documents present the same content or topic, so that SW4ALL can be a rich resource for aiding teachers and learners that wish to compare different linguistic strategies for writing a similar content, providing an interesting option for improving the learning of English as a second language.

## 7. Acknowledgements

The authors would like to thank the Walloon Region (Projects BEWARE n. 1510637 and 1610378) for support, and Altissia International for research collaboration.

## 8. Bibliographical References

Arthur, D. and Vassilvitskii, S. (2007). k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035. Society for Industrial and Applied Mathematics.

Bruce, B., Rubin, A., and Starr, K. (1981). Why readability formulas fail. *IEEE Transactions on Professional Communication*, (1):50–52.

Callan, J. and Eskenazi, M. (2007). Combining lexical and grammatical features to improve readability measures for first and second language texts. In *Proceedings of NAACL HLT*, pages 460–467. sn.

Chinkina, M. and Meurers, D. (2016). Linguistically aware information retrieval: Providing input enrichment for second language learners. In *BEA@ NAACL-HLT*, pages 188–198.

Chinkina, M., Kannan, M., and Meurers, D. (2016). Online information retrieval for language learning. *ACL 2016*, page 7.

Coster, W. and Kauchak, D. (2011). Learning to simplify sentences using wikipedia. In *Proceedings of the workshop on monolingual text-to-text generation*, pages 1–9. Association for Computational Linguistics.

François, T. and Fairon, C. (2012). An ai readability formula for french as a foreign language. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 466–477. Association for Computational Linguistics.

Frank, E., Hall, M., Holmes, G., Kirkby, R., Pfahringer, B., Witten, I. H., and Trigg, L. (2009). Weka-a machine learning workbench for data mining. In *Data mining and knowledge discovery handbook*, pages 1269–1277. Springer.

Gunning, R. (1952). The technique of clear writing.

Heilman, M., Zhao, L., Pino, J., and Eskenazi, M. (2008). Retrieval of reading materials for vocabulary and reading practice. In *Proceedings of the Third Workshop on Innovative Use of NLP for Building Educational Applications*, pages 80–88. Association for Computational Linguistics.

Kauchak, D. (2013). Improving text simplification language modeling using unsimplified text data. In *ACL (1)*, pages 1537–1546.

Kincaid, J. P., Fishburne Jr, R. P., Rogers, R. L., and Chissom, B. S. (1975). Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. Technical report, Naval Technical Training Command Millington TN Research Branch.

Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J. R., Bethard, S., and McClosky, D. (2014). The stanford corenlp natural language processing toolkit. In *ACL (System Demonstrations)*, pages 55–60.

Miltsakaki, E. and Trount, A. (2007). Read-x: Automatic evaluation of reading difficulty of web text. In *E-Learn: World Conference on E-Learning in Corporate, Government, Healthcare, and Higher Education*, pages 7280–7286. Association for the Advancement of Computing in Education (AACE).

Ott, N. and Meurers, D. (2011). Information retrieval for education: Making search engines language aware. *Themes in Science and Technology Education*, 3(1-2):9–30.

Passonneau, R., Hemat, L., Plante, J., and Sheehan, K. M. (2002). Electronic sources as input to gre® reading

<sup>9</sup>The k value was set as 10.

- comprehension item development: Sourcefinder prototype evaluation. *ETS Research Report Series*, 2002(1).
- Pilán, I., Volodina, E., and Zesch, T. (2016). Predicting proficiency levels in learner writings by transferring a linguistic complexity model from expert-written coursebooks. In *COLING*, pages 2101–2111.
- Rott, S. (1999). The effect of exposure frequency on intermediate language learners incidental vocabulary acquisition and retention through reading. *Studies in second language acquisition*, 21(4):589–619.
- Sheehan, K. M., Kostin, I., and Futagi, Y. (2007). Sourcefinder: A construct-driven approach for locating appropriately targeted reading comprehension source texts. In *Workshop on Speech and Language Technology in Education*.
- Vajjala, S. and Meurers, D. (2013). On the applicability of readability models to web texts. In *Proceedings of the 2nd Workshop on Predicting and Improving Text Readability for Target Reader Populations*, pages 59–68.
- Verhelst, N., Van Avermaet, P., Takala, S., Figueras, N., and North, B. (2009). *Common European Framework of Reference for Languages: learning, teaching, assessment*. Cambridge University Press.
- Wisniewski, K., Schöne, K., Nicolas, L., Vettori, C., Boyd, A., Meurers, D., Abel, A., and Hana, J. (2013). Merlin: An online trilingual learner corpus empirically grounding the european reference levels in authentic learner data. In *ICT for Language Learning 2013, Conference Proceedings, Florence, Italy. Libreriauniversitaria. it Edizioni*.
- Xia, M., Kochmar, E., and Briscoe, E. (2016). Text readability assessment for second language learners.
- Zilio, L. and Fairon, C. (2017). Adaptive system for language learning. In *Advanced Learning Technologies (ICALT), 2017 IEEE 17th International Conference on*, pages 47–49. IEEE.
- Zilio, L., Wilkens, R., and Fairon, C. (2017a). Enhancing grammatical structures in web-based texts. In *Proceedings of the 25th EUROCALL*, pages 839–846. Accepted.
- Zilio, L., Wilkens, R., and Fairon, C. (2017b). Using nlp for enhancing second language acquisition. In *Proceedings of Recent Advances in Natural Language Processing*, pages 839–846.

## 9. Language Resource References

- Geertzen, J., Alexopoulou, T., and Korhonen, A. (2013). Automatic linguistic annotation of large scale 12 databases: The ef-cambridge open language database (efcamdat). In *Proceedings of the 31st Second Language Research Forum. Somerville, MA: Cascadilla Proceedings Project*.