

Enriching Complex Networks with Word Embeddings for Detecting Mild Cognitive Impairment from Speech Transcripts

Leandro B. dos Santos¹, Edilson A. Corrêa Jr¹, Osvaldo N. Oliveira Jr², Diego R. Amancio¹,
Letícia L. Mansur³, Sandra M. Aluísio¹

¹ Institute of Mathematics and Computer Science, University of São Paulo, São Carlos, São Paulo, Brazil

² São Carlos Institute of Physics, University of São Paulo, São Carlos, São Paulo, Brazil

³ Department of Physiotherapy, Speech Pathology and Occupational Therapy,
University of São Paulo, São Paulo, São Paulo, Brazil

{leandrob, edilsonacjr, lamansur}@usp.br, chu@ifsc.usp.br
{diego, sandra}@icmc.usp.br

Abstract

Mild Cognitive Impairment (MCI) is a mental disorder difficult to diagnose. Linguistic features, mainly from parsers, have been used to detect MCI, but this is not suitable for large-scale assessments. MCI disfluencies produce non-grammatical speech that requires manual or high precision automatic correction of transcripts. In this paper, we modeled transcripts into complex networks and enriched them with word embedding (CNE) to better represent short texts produced in neuropsychological assessments. The network measurements were applied with well-known classifiers to automatically identify MCI in transcripts, in a binary classification task. A comparison was made with the performance of traditional approaches using Bag of Words (BoW) and linguistic features for three datasets: DementiaBank in English, and Cinderella and Arizona-Battery in Portuguese. Overall, CNE provided higher accuracy than using only complex networks, while Support Vector Machine was superior to other classifiers. CNE provided the highest accuracies for DementiaBank and Cinderella, but BoW was more efficient for the Arizona-Battery dataset probably owing to its short narratives. The approach using linguistic features yielded higher accuracy if the transcriptions of the Cinderella dataset were manually revised. Taken together, the results indicate that complex networks enriched with embedding is promising for detecting MCI in large-scale assessments.

1 Introduction

Mild Cognitive Impairment (MCI) can affect one or multiple cognitive domains (e.g. memory, language, visuospatial skills and executive functions), and may represent a pre-clinical stage of Alzheimer's disease (AD). The impairment that affects memory, referred to as amnesic MCI, is the most frequent, with the highest conversion rate for AD, at 15% per year versus 1 to 2% for the general population. Since dementias are chronic and progressive diseases, their early diagnosis ensures a greater chance of success to engage patients in non-pharmacological treatment strategies such as cognitive training, physical activity and socialization (Teixeira et al., 2012).

Language is one of the most efficient information sources to assess cognitive functions. Changes in language usage are frequent in patients with dementia and are normally first recognized by the patients themselves or their family members. Therefore, the automatic analysis of discourse production is promising in diagnosing MCI at early stages, which may address potentially reversible factors (Muangpaisan et al., 2012). Proposals to detect language-related impairment in dementias include machine learning (Jarrold et al., 2010; Roark et al., 2011; Fraser et al., 2014, 2015), magnetic resonance imaging (Dyrba et al., 2015), and data screening tests added to demographic information (Weakley et al., 2015). Discourse production (mainly narratives) is attractive because it allows the analysis of linguistic microstructures, including phonetic-phonological, morphosyntactic and semantic-lexical components, as well as semantic-pragmatic macrostructures.

Automated discourse analysis based on Natural Language Processing (NLP) resources and tools to diagnose dementias via machine learning methods has been used for English language (Lehr et al.,

2012; Jarrold et al., 2014; Orimaye et al., 2014; Fraser et al., 2015; Davy et al., 2016) and for Brazilian Portuguese (Aluísio et al., 2016). A variety of features are required for this analysis, including Part-of-Speech (PoS), syntactic complexity, lexical diversity and acoustic features. Producing robust tools to extract these features is extremely difficult because speech transcripts used in neuropsychological evaluations contain disfluencies (repetitions, revisions, paraphasias) and patient’s comments about the task being evaluated. Another problem in using linguistic knowledge is the high dependence on manually created resources, such as hand-crafted linguistic rules and/or annotated corpora. Even when traditional statistical techniques (Bag of Words or ngrams) are applied, problems still appear in dealing with disfluencies, because mispronounced words will not be counted together. Indeed, other types of disfluencies (repetition, amendments, patient’s comments about the task) will be counted, thus increasing the vocabulary.

An approach applied successfully to several areas of NLP (Mihalcea and Radev, 2011), which may suffer less from the problems mentioned above, relies on the use of complex networks and graph theory. The word adjacency network model (i Cancho and Solé, 2001; Roxas and Tapang, 2010; Amancio et al., 2012a; Amancio, 2015b) has provided good results in text classification (de Arruda et al., 2016) and related tasks, namely author detection (Amancio, 2015a), identification of literary movements (Amancio et al., 2012c), authenticity verification (Amancio et al., 2013) and word sense discrimination (Amancio et al., 2012b).

In this paper, we show that speech transcripts (narratives or descriptions) can be modeled into complex networks that are enriched with word embedding in order to better represent short texts produced in these assessments. When applied to a machine learning classifier, the complex network features were able to distinguish between control participants and mild cognitive impairment participants. Discrimination of the two classes could be improved by combining complex networks with linguistic and traditional statistical features.

With regard to the task of detecting MCI from transcripts, this paper is, to the best of our knowledge, the first to: a) show that classifiers using features extracted from transcripts modeled into

complex networks enriched with word embedding present higher accuracy than using only complex networks for 3 datasets; and b) show that for languages that do not have competitive dependency and constituency parsers to exploit syntactic features, e.g. Brazilian Portuguese, complex networks enriched with word embedding constitute a source to extract new, language independent features from transcripts.

2 Related Work

Detection of memory impairment has been based on linguistic, acoustic, and demographic features, in addition to scores of neuropsychological tests. Linguistic and acoustic features were used to automatically detect aphasia (Fraser et al., 2014); and AD (Fraser et al., 2015) or dementia (Orimaye et al., 2014) in the public corpora of Dementia-Bank¹. Other studies distinguished different types of dementia (Garrard et al., 2014; Jarrold et al., 2014), in which speech samples were elicited using the Picnic picture of the Western Aphasia Battery (Kertesz, 1982). Davy et al. (2016) also used the Picnic scene to detect MCI, where the subjects were asked to write (by hand) a detailed description of the scene.

As for automatic detection of MCI in narrative speech, Roark et al. (2011) extracted speech features and linguistic complexity measures of speech samples obtained with the Wechsler Logical Memory (WLM) subtest (Wechsler et al., 1997), and Lehr et al. (2012) fully automatized the WLM subtest. In this test, the examiner tells a short narrative to a subject, who then retells the story to the examiner, immediately and after a 30-minute delay. WLM scores are obtained by counting the number of story elements recalled.

Tóth et al. (2015) and Vincze et al. (2016) used short animated films to evaluate immediate and delayed recalls in MCI patients who were asked to talk about the first film shown, then about their previous day, and finally about another film shown last. Tóth et al. (2015) adopted automatic speech recognition (ASR) to extract a phonetic level segmentation, which was used to calculate acoustic features. Vincze et al. (2016) used speech, morphological, semantic, and demographic features collected from their speech transcripts to automatically identify patients suffering from MCI.

For the Portuguese language, machine learning

¹talkbank.org/DementiaBank/

algorithms were used to identify subjects with AD and MCI. [Aluísio et al. \(2016\)](#) used a variety of linguistic metrics, such as syntactic complexity, idea density ([da Cunha et al., 2015](#)), and text cohesion through latent semantics. NLP tools with high precision are needed to compute these metrics, which is a problem for Portuguese since no robust dependency or constituency parsers exist. Therefore, the transcriptions had to be manually revised; they were segmented into sentences, following a semantic-structural criterion and capitalization was applied. The authors also removed disfluencies and inserted omitted subjects when they were hidden, in order to reduce parsing errors. This process is obviously expensive, which has motivated us to use complex networks in the present study to model transcriptions and avoid a manual preprocessing step.

3 Modeling and Characterizing Texts as Complex Networks

The theory and concepts of complex networks have been used in several NLP tasks ([Mihalcea and Radev, 2011](#); [Cong and Liu, 2014](#)), such as text classification ([de Arruda et al., 2016](#)), summarization ([Antiqueira et al., 2009](#); [Amancio et al., 2012a](#)) and word sense disambiguation ([Silva and Amancio, 2012](#)). In this study, we used the word co-occurrence model (also called word adjacency model) because most of the syntactical relations occur among neighboring words ([i Cancho et al., 2004](#)). Each distinct word becomes a node and words that are adjacent in the text are connected by an edge. Mathematically, a network is defined as an undirected graph $G = \{V, E\}$, formed by a set $V = \{v_1, v_2, \dots, v_n\}$ of nodes (words) and a set $E = \{e_1, e_2, \dots, e_m\}$ of edges (co-occurrence) that are represented by an adjacency matrix A , whose elements A_{ij} are equal to 1 whenever there is an edge connecting nodes (words) i and j , and equal to 0 otherwise.

Before modeling texts into complex networks, it is often necessary to do some preprocessing in the raw text. Preprocessing starts with tokenization where each document/text is divided into tokens (meaningful elements, e.g., words and punctuation marks) and then *stopwords* and punctuation marks are removed, since they have little semantic meaning. One last step we decided to eliminate from the preprocessing pipeline is lemmatization, which transforms each word into its canonical

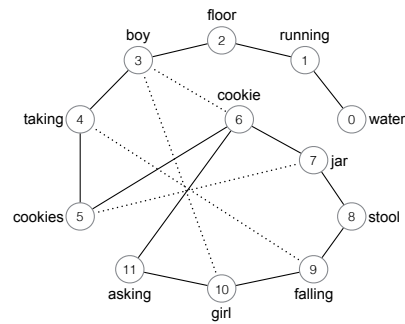


Figure 1: Example of co-occurrence network enriched with semantic information for the following transcription: “*The water’s running on the floor. Boy’s taking cookies out of cookie out of the cookie jar. The stool is falling over. The girl was asking for a cookie.*”. The solid edges of the network represent co-occurrence edges and the dotted edges represent connections between words that had similarity higher than 0.5.

form. This decision was made based on two factors. First, a recent work has shown that lemmatization has little or no influence when network modeling is adopted in related tasks ([Machicao et al., 2016](#)). Second, the lemmatization process requires part-of-speech (POS) tagging that may introduce undesirable noises/errors in the text, since the transcriptions in our work contain disfluencies.

Another problem with transcriptions in our work is their size. As demonstrated by [Amancio \(2015c\)](#), classification of small texts using networks can be impaired, since short texts have almost linear networks, and the topological measures of these networks have little or no information relevant to classification. To solve this problem, we adapted the approach of inducing language networks from word embeddings, proposed by [Perozzi et al. \(2014\)](#) to enrich the networks with semantic information. In their work, language networks were generated from continuous word representations, in which each word is represented by a dense, real-valued vector obtained by training neural networks in the language model task (or variations, such as context prediction) ([Benigno et al., 2003](#); [Collobert et al., 2011](#); [Mikolov et al., 2013a,b](#)). This structure is known to capture syntactic and semantic information. [Perozzi et al. \(2014\)](#), in particular, take advantage of word embeddings to build networks where each word is

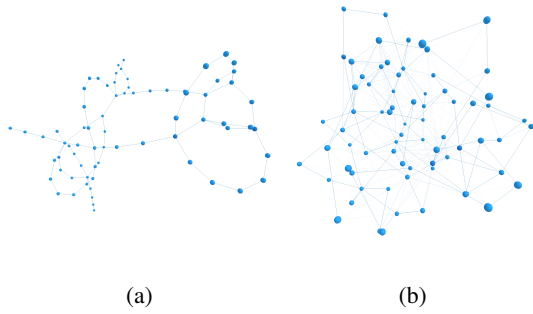


Figure 2: Example of (a) co-occurrence network created for a transcript of the Cookie Theft dataset (see Supplementary Information, Section A) and (b) the same co-occurrence network enriched with semantic information. Note that (b) is a more informative network than (a), since (a) is practically a linear network.

a vertex and edges are defined by similarity between words established by the proximity of the word vectors.

Following this methodology, in our model we added new edges to the co-occurrence networks considering similarities between words, that is, for all pairs of words in the text that were not connected, an edge was created if their vectors (from word embedding) had a cosine similarity higher than a given threshold. Figure 1 shows an example of a co-occurrence network enriched by similarity links (the dotted edges). The gain in information by enriching a co-occurrence network with semantic information is readily apparent in Figure 2.

4 Datasets, Features and Methods

4.1 Datasets

The datasets² used in our study consisted of: (i) manually segmented and transcribed samples from the DementiaBank and Cinderella story and (ii) transcribed samples of Arizona Battery for Communication Disorders of Dementia (ABCD) automatically segmented into sentences, since we are working towards a fully automated system to detect MCI in transcripts and would like to evaluate a dataset which was automatically processed.

The DementiaBank dataset is composed of short English descriptions, while the Cinderella dataset contains longer Brazilian Portuguese narratives. ABCD dataset is composed of very short narratives, also in Portuguese. Below, we describe

²All datasets are made available in the same representations used in this work, upon request to the authors.

in further detail the datasets, participants, and the task in which they were used.

4.1.1 The Cookie Theft Picture Description Dataset

The clinical dataset used for the English language was created during a longitudinal study conducted by the University of Pittsburgh School of Medicine on Alzheimer’s and related dementia, funded by the National Institute of Aging. To be eligible for inclusion in the study, all participants were required to be above 44 years of age, have at least 7 years of education, no history of nervous system disorders nor be taking neuroleptic medication, have an initial Mini-Mental State Exam (MMSE) score of 10 or greater, and be able to give informed consent. The dataset contains transcripts of verbal interviews with AD and related Dementia patients, including those with MCI (for further details see (Becker et al., 1994)).

We used 43 transcriptions with MCI in addition to another 43 transcriptions sampled from 242 healthy elderly people to be used as the control group. Table 1 shows the demographic information for the two diagnostic groups.

Demographic	Control	MCI
Avg. Age (SD)	64.1 (7.2)	69.3 (8.2)
No. of Male/Female	23/20	27/16

Table 1: Demographic information of participants in the Cookie Theft dataset.

For this dataset, interviews were conducted in English and narrative speech was elicited using the Cookie Theft picture (Goodglass et al., 2001) (Figure 3 from Goodglass et al. (2001) in Section A.1). During the interview, patients were given the picture and were told to discuss everything they could see happening in the picture. The patients’ verbal utterances were recorded and then transcribed into the CHAT (Codes for the Human Analysis of Transcripts) transcription format (MacWhinney, 2000).

We extracted the word-level transcript patient sentences from the CHAT files and discarded the annotations, as our goal was to create a fully automated system that does not require the input of a human annotator. We automatically removed filled pauses such as *uh*, *um*, *er*, and *ah* (e.g. *uh it seems to be summer out*), short false starts (e.g. *just t the ones*), and repetition (e.g. *mother’s finished certain of the the dishes*), as in (Fraser et al.,

2015). The control group had an average of 9.58 sentences per narrative, with each sentence having an average of 9.18 words; while the MCI group had an average of 10.97 sentences per narrative, with 10.33 words per sentence in average.

4.1.2 The Cinderella Narrative Dataset

The dataset examined in this study included 20 subjects with MCI and 20 normal elderly control subjects, as diagnosed at the Medical School of the University of São Paulo (FMUSP). Table 2 shows the demographic information of the two diagnostic groups, which were also used in Aluísio et al. (2016).

Demographic	Control	MCI
Avg. Age (SD)	74.8 (11.3)	73.3 (5.9)
Avg. Years of Education (SD)	11.4 (2.6)	10.8 (4.5)
No. of Male/Female	27/16	29/14

Table 2: Demographic information of participants in the Cinderella dataset.

The criteria used to diagnose MCI came from Petersen (2004). Diagnostics were carried out by a multidisciplinary team consisting of psychiatrists, geriatricians, neurologists, neuropsychologists, speech pathologists, and occupational therapists, by a criterion of consensus. Inclusion criteria for the control group were elderly with no cognitive deficits and preservation of functional capacity in everyday life. The exclusion criteria for the normal group were: poorly controlled clinical diseases, sensitive deficits that were not being compensated for and interfered with the performance in tests, and other neurological or psychiatric diagnoses associated with dementia or cognitive deficits and use of medications in doses that affected cognition.

Speech narrative samples were elicited by having participants tell the Cinderella story; participants were given as much time as they needed to examine a picture book illustrating the story (Figure 4 in Section A). When each participant had finished looking at the pictures, the examiner asked the subject to tell the story in their own words, as in Saffran et al. (1989). The time was recorded, but there was no limit imposed to the narrative length. If the participant had difficulty initiating or continuing speech, or took a long pause, an evaluator would use the stimulus question “What happens next?”, seeking to encourage the participant to continue his/her narrative. When the sub-

ject was unable to proceed with the narrative, the examiner asked if he/she had finished the story and had something to add. Each speech sample was recorded and then manually transcribed at the word level following the NURC/SP N. 338 EF and 331 D2 transcription norms³.

Other tests were applied after the narrative, in the following sequence: phonemic verbal fluency test, action verbal fluency, Camel and Cactus test (Bozeat et al., 2000), and Boston Naming test (Kaplan et al., 2001), in order to diagnose the groups.

Since our ultimate goal is to create a fully automated system that does not require the input of a human annotator, we manually segmented sentences to simulate a high-quality ASR transcript with sentence segmentation, and we automatically removed the disfluencies following the same guidelines of TalkBank project. However, other disfluencies (revisions, elaboration, paraphasias and comments about the task) were kept. The control group had an average of 30.80 sentences per narrative, and each sentence averaged 12.17 words. As for the MCI group, it had an average of 29.90 sentences per narrative, and each sentence averaged 13.03 words.

We also evaluated a different version of the dataset used in Aluísio et al. (2016), where narratives were manually annotated and revised to improve parsing results. The revision process was the following: (i) in the original transcript, segments with hesitations or repetitions of more than one word or segment of a single word were annotated to become a feature and then removed from the narrative to allow the extraction of features from parsing; (ii) empty emissions, which were comments unrelated to the topic of narration or confirmations, such as “né” (alright), were also annotated and removed; (iii) prolongations of vowels, short pauses and long pauses were also annotated and removed; and (iv) omitted subjects in sentences were inserted. In this revised dataset, the control group had an average of 45.10 sentences per narrative, and each sentence averaged 8.17 words. The MCI group had an average of 31.40 sentences per narrative, with each sentence averaging 10.91 words.

4.1.3 The ABCD Dataset

The subtest of immediate/delayed recall of narratives of the ABCD battery was administered to 23

³albertofedel.blogspot.com.br/2010_11_01_archive.html

participants with a diagnosis of MCI and 20 normal elderly control participants, as diagnosed at the Medical School of the University of São Paulo (FMUSP).

MCI subjects produced 46 narratives while the control group produced 39 ones. In order to carry out experiments with a balanced corpus, as with the previous two datasets, we excluded seven transcriptions from the MCI group. We used the automatic sentence segmentation method referred to as DeepBond (Treviso et al., 2017) in the transcripts.

Table 3 shows the demographic information. The control group had an average of 5.23 sentences per narrative, with 11 words per sentence on average, and the MCI group had an average of 4.95 sentences per narrative, with an average of 12.04 words per sentence. Interviews were conducted in Portuguese and the subject listened to the examiner read a short narrative. The subject then retold the narrative to the examiner twice: once immediately upon hearing it and again after a 30-minute delay (Bayles and Tomoeda, 1991). Each speech sample was recorded and then manually transcribed at the word level following the NURC/SP N. 338 EF and 331 D2 transcription norms.

Demographic	Control	MCI
Avg. Age (SD)	61 (7.5)	72,0 (7.4)
Avg. Years of Education (SD)	16 (7.6)	13.3 (4.2)
No. of Male/Female	6/14	16/7

Table 3: Demographic information of participants in the ABCD dataset.

4.2 Features

Features of three distinct natures were used to classify the transcribed texts: topological metrics of co-occurrence networks, linguistic features and bag of words representations.

4.2.1 Topological Characterization of Networks

Each transcription was mapped into a co-occurrence network, and then enriched via word embeddings using the cosine similarity of words. Since the occurrence of out-of-vocabulary words is common in texts of neuropsychological assessments, we used the method proposed by Bojanowski et al. (2016) to generate word embeddings. This method extends the skip-gram model to use character-level information, with each word

being represented as a bag of character n -grams. It provides some improvement in comparison with the traditional skip-gram model in terms of syntactic evaluation (Mikolov et al., 2013b) but not for semantic evaluation.

Once the network has been enriched, we characterize its topology using the following ten measurements:

1. **PageRank:** is a centrality measurement that reflects the relevance of a node based on its connections to other relevant nodes (Brin and Page, 1998);
2. **Betweenness:** is a centrality measurement that considers a node as relevant if it is highly accessed via shortest paths. The betweenness of a node v is defined as the fraction of shortest paths going through node v ;
3. **Eccentricity:** of a node is calculated by measuring the shortest distance from the node to all other vertices in the graph and taking the maximum;
4. **Eigenvector centrality:** is a measurement that defines the importance of a node based on its connectivity to high-rank nodes;
5. **Average Degree of the Neighbors of a Node:** is the average of the degrees of all its direct neighbors;
6. **Average Shortest Path Length of a Node:** is the average distance between this node and all other nodes of the network;
7. **Degree:** is the number of edges connected to the node;
8. **Assortativity Degree:** or degree correlation measures the tendency of nodes to connect to other nodes that have similar degree;
9. **Diameter:** is defined as the maximum shortest path;
10. **Clustering Coefficient:** measures the probability that two neighbors of a node are connected.

Most of the measurements described above are local measurements, i.e. each node i possesses a value X_i , so we calculated the average $\mu(X)$, standard deviation $\sigma(X)$ and skewness $\gamma(X)$ for each measurement (Amancio, 2015b).

4.2.2 Linguistic Features

Linguistic features for classification of neuropsychological assessments have been used in several studies (Roark et al., 2011; Jarrold et al., 2014; Fraser et al., 2014; Orimaye et al., 2014; Fraser et al., 2015; Vincze et al., 2016; Davy et al., 2016). We used the Coh-Metrix⁴(Graesser et al., 2004) tool to extract features from English transcripts, resulting in 106 features. The metrics are divided into eleven categories: Descriptive, Text Easability Principal Component, Referential Cohesion, Latent Semantic Analysis (LSA), Lexical Diversity, Connectives, Situation Model, Syntactic Complexity, Syntactic Pattern Density, Word Information, and Readability (Flesch Reading Ease, Flesch-Kincaid Grade Level, Coh-Metrix L2 Readability).

For Portuguese, Coh-Metrix-Dementia (Aluísio et al., 2016) was used. The metrics affected by constituency and dependency parsing were not used because they are not robust with disfluencies. Metrics based on manual annotation (such as proportion short pauses, mean pause duration, mean number of empty words, and others) were also discarded. The metrics of Coh-Metrix-Dementia are divided into twelve categories: Ambiguity, Anaphoras, Basic Counts, Connectives, Co-reference Measures, Content Word Frequencies, Hypernyms, Logic Operators, Latent Semantic Analysis, Semantic Density, Syntactical Complexity, and Tokens. The metrics used are shown in detail in Section A.2. In total, 58 metrics were used, from the 73 available on the website⁵.

4.2.3 Bag of Words

The representation of text collections under the BoW assumption (i.e., with no information relating to word order) has been a robust solution for text classification. In this methodology, transcripts are represented by a table in which the columns represent the terms (or existing words) in the transcripts and the values represent frequency of a term in a document.

4.3 Classification Algorithms

In order to quantify the ability of the topological characterization of networks, linguistic metrics and BoW features were used to distinguish subjects with MCI from healthy controls. We

employed four machine learning algorithms to induce classifiers from a training set. These techniques were the Gaussian Naive Bayes (G-NB), k -Nearest Neighbor (k -NN), Support Vector Machine (SVM), linear and radial bases functions (RBF), and Random Forest (RF). We also combined these classifiers through ensemble and multi-view learning. In ensemble learning, multiple models/classifiers are generated and combined using a majority vote or the average of class probabilities to produce a single result (Zhou, 2012).

In multi-view learning, multiple classifiers are trained in different feature spaces and thus combined to produce a single result. This approach is an elegant solution in comparison to combining all features in the same vector or space, for two main reasons. First, combination is not a straightforward step and may lead to noise insertion since the data have different natures. Second, using different classifiers for each feature space allows for different weights to be given for each type of feature, and these weights can be learned by a regression method to improve the model. In this work, we used majority voting to combine different feature spaces.

5 Experiments and Results

All experiments were conducted using the Scikit-learn⁶ (Pedregosa et al., 2011), with classifiers evaluated on the basis of classification accuracy i.e. the total proportion of narratives which were correctly classified. The evaluation was performed using 5-fold cross-validation instead of the well-accepted 10-fold cross-validation because the datasets in our study were small and the test set would have shrunk, leading to less precise measurements of accuracy. The threshold parameter was optimized with the best values being 0.7 in the Cookie Theft dataset and 0.4 in both the Cinderella and ABCD datasets.

We used the model proposed by Bojanowski et al. (2016) with default parameters (100 dimensional embeddings, context window equal to 5 and 5 epochs) to generate word embedding. We trained the models in Portuguese and English Wikipedia dumps from October and November 2016 respectively.

The accuracy in classification is given in Tables 4 through 6. CN, CNE, LM, and BoW denote, respectively, complex networks, complex network

⁴cohmetrix.com

⁵<http://143.107.183.175:22380>

⁶<http://scikit-learn.org>

enriched with embedding, linguistic metrics and Bag of Words, and CNE-LM, CNE-BoW, LM-BoW and CNE-LM-BoW refer to combinations of the feature spaces (multiview learning), using the majority vote. Cells with the “-” sign mean that it was not possible to apply majority voting because there were two classifiers. The last line represents the use of an ensemble of machine learning algorithms, in which the combination used was the majority voting in both ensemble and multiview learning.

In general, CNE outperforms the approach using only complex networks (CN), while SVM (Linear or RBF kernel) provides higher accuracy than other machine learning algorithms. The results for the three datasets show that characterizing transcriptions into complex networks is competitive with other traditional methods, such as the use of linguistic metrics. In fact, among the three types of features, using enriched networks (CNE) provided the highest accuracies in two datasets (Cookie Theft and original Cinderella). For the ABCD dataset, which contains short narratives, the small length of the transcriptions may have had an effect, since BoW features led to the highest accuracy. In the case of the revised Cinderella dataset, segmented into sentences and capitalized as reported in [Aluísio et al. \(2016\)](#), Table 7 shows that the manual revision was an important factor, since the highest accuracies were obtained with the approach based on linguistic metrics (LM). However, this process of manually removing disfluencies demands time; therefore it is not practical for large-scale assessments.

Ensemble and multi-view learning were helpful for the Cookie Theft dataset, in which multi-view learning achieved the highest accuracy (65% of accuracy for narrative texts, a 3% of improvement compared to the best individual classifier). However, neither multi-view or ensemble learning enhanced accuracy in the Cinderella dataset, where SVM-RBF with CNE space achieved the highest accuracy (65%). For the ABCD dataset, multi-view CNE-LM-BoW with SVM-RBF and KNN classifiers improved the accuracy to 4% and 2%, respectively. Somewhat surprising were the results of SVM with linear kernel in BoW feature space (75% of accuracy).

6 Conclusions and Future Work

In this study, we employed metrics of topological properties of CN in a machine learning classification approach to distinguish between healthy patients and patients with MCI. To the best of our knowledge, these metrics have never been used to detect MCI in speech transcripts; CN were enriched with word embeddings to better represent short texts produced in neuropsychological assessments. The topological properties of CN outperform traditional linguistic metrics in individual classifiers’ results. Linguistic features depend on grammatical texts to present good results, as can be seen in the results of the manually processed Cinderella dataset (Table 7). Furthermore, we found that combining machine and multi-view learning can improve accuracy. The accuracies found here are comparable to the values reported by other authors, ranging from 60% to 85% ([Prud’hommeaux and Roark, 2011](#); [Lehr et al., 2012](#); [Tóth et al., 2015](#); [Vincze et al., 2016](#)), which means that it is not easy to distinguish between healthy subjects and those with cognitive impairments. The comparison with our results is not straightforward, though, because the databases used in the studies are different. There is a clear need for publicly available datasets to compare different methods, which would optimize the detection of MCI in elderly people.

In future work, we intend to explore other methods to enrich CN, such as the Recurrent Language Model, and use other metrics to characterize an adjacency network. The pursuit of these strategies is relevant because language is one of the most efficient information sources to evaluate cognitive functions, commonly used in neuropsychological assessments. As this work is ongoing, we will keep collecting new transcriptions of the ABCD retelling subtest to increase the corpus size and obtain more reliable results in our studies. Our final goal is to apply neuropsychological assessment batteries, such as the ABCD retelling subtest, to mobile devices, specifically tablets. This adaptation will enable large-scale applications in hospitals and facilitate the maintenance of application history in longitudinal studies, by storing the results in databases immediately after the test application.

Classifier	CN	CNE	LM	BoW	CNE-LM	CNE-BoW	LM-BoW	CNE-LM-BoW
SVM-Linear	52	55	56	59	–	–	–	60
SVM-RBF	56	62	58	60	–	–	–	65
<i>k</i> -NN	59	61	46	57	–	–	–	59
RF	52	47	45	48	–	–	–	50
G-NB	51	48	56	55	–	–	–	50
Ensemble	56	60	54	58	57	60	63	65

Table 4: Classification accuracy achieved on Cookie Theft dataset.

Classifier	CN	CNE	LM	BoW	CNE-LM	CNE-BoW	LM-BoW	CNE-LM-BoW
SVM-Linear	52	60	52	50	–	–	–	52
SVM-RBF	57	65	47	37	–	–	–	50
<i>k</i> -NN	47	50	47	37	–	–	–	37
RF	55	57	47	45	–	–	–	52
G-NB	47	52	47	55	–	–	–	52
Ensemble	52	60	50	37	57	52	50	47

Table 5: Classification accuracy achieved on Cinderella dataset.

Classifier	CN	CNE	LM	BoW	CNE-LM	CNE-BoW	LM-BoW	CNE-LM-BoW
SVM-Linear	56	69	51	75	–	–	–	74
SVM-RBF	54	57	66	67	–	–	–	71
<i>k</i> -NN	56	56	69	63	–	–	–	71
RF	54	62	70	64	–	–	–	69
G-NB	61	55	55	65	–	–	–	65
Ensemble	55	61	62	72	69	68	75	73

Table 6: Classification accuracy achieved on ABCD dataset.

Classifier	CN	CNE	LM	BoW
SVM-Linear	50	65	65	52
SVM-RBF	57	67	72	55
KNN	42	47	55	50
RF	52	47	70	45
G-NB	52	65	62	45
Ensemble	52	60	72	45

Table 7: Classification accuracy achieved on Cinderella dataset manually processed to revise non-grammatical sentences.

Acknowledgments

This work was supported by CAPES, CNPq, FAPESP, and Google Research Awards in Latin America. We would like to thank NVIDIA for their donation of GPU.

References

Sandra M. Aluísio, Andre L. da Cunha, and Carolina Scarton. 2016. [Evaluating progression of alzheimer’s disease by regression and classification methods in a narrative language test in portuguese](#). In João Silva, Ricardo Ribeiro, Paulo Quaresma, André Adami, and António Branco, editors, *Internation*

ational Conference on Computational Processing of the Portuguese Language. Springer, pages 109–114. https://doi.org/10.1007/978-3-319-41552-9_10.

Diego R. Amancio. 2015a. [Authorship recognition via fluctuation analysis of network topology and word intermittency](#). *Journal of Statistical Mechanics: Theory and Experiment* 2015(3):P03005. <https://doi.org/10.1088/1742-5468/2015/03/P03005>.

Diego R. Amancio. 2015b. [A complex network approach to stylometry](#). *PloS one* 10(8):e0136076. <https://doi.org/10.1371/journal.pone.0136076>.

Diego R. Amancio. 2015c. [Probing the topological properties of complex networks modeling short written texts](#). *PloS one* 10(2):1–17. <https://doi.org/10.1371/journal.pone.0118394>.

Diego R. Amancio, Eduardo G. Altmann, Diego Rybski, Osvaldo N. Oliveira Jr., and Luciano da F. Costa. 2013. [Probing the statistical properties of unknown texts: Application to the voynich manuscript](#). *PLOS ONE* 8(7):1–10. <https://doi.org/10.1371/journal.pone.0067310>.

Diego R. Amancio, Maria G. V. Nunes, Osvaldo N. Oliveira Jr., and Luciano F. Costa. 2012a. [Extractive summarization using complex networks and syntactic dependency](#). *Physica A: Statistical Me-*

- chanics and its Applications* 391(4):1855–1864. <https://doi.org/10.1016/j.physa.2011.10.015>.
- Diego R. Amancio, O.N. Oliveira Jr., and Luciano da F. Costa. 2012b. Unveiling the relationship between complex networks metrics and word senses. *EPL (Europhysics Letters)* 98(1):18002. <https://doi.org/10.1209/0295-5075/98/18002>.
- Diego R. Amancio, Osvaldo N. Oliveira Jr., and Luciano F. Costa. 2012c. Identification of literary movements using complex networks to represent texts. *New Journal of Physics* 14(4):043029. <https://doi.org/10.1088/1367-2630/14/4/043029>.
- Lucas Antiquiera, Osvaldo N. Oliveira Jr., Luciano da Fontoura Costa, and Maria das Graças Volpe Nunes. 2009. A complex network approach to text summarization. *Information Sciences* 179(5):584 – 599.
- Kathryn A. Bayles and Cheryl K. Tomoeda. 1991. *ABCD: Arizona Battery for Communication Disorders of Dementia*. Tucson, AZ: Canyonlands Publishing.
- James T. Becker, François Boiler, Oscar L. Lopez, Judith Saxton, and Karen L. McGonigle. 1994. The natural history of alzheimer’s disease: description of study cohort and accuracy of diagnosis. *Archives of Neurology* 51(6):585–594. <https://doi.org/10.1001/archneur.1994.00540180063015>.
- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *journal of machine learning research* 3(Feb):1137–1155.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- Sasha Bozeat, Matthew A. Ralph, Karalyn Patterson, Peter Garrard, and John R. Hodges. 2000. Non-verbal semantic impairment in semantic dementia. *Neuropsychologia* 38(9):1207–1215. [https://doi.org/10.1016/S0028-3932\(00\)00034-8](https://doi.org/10.1016/S0028-3932(00)00034-8).
- Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual web search engine. In *International Conference on World Wide Web*. Elsevier, pages 107–117.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research* 12(Aug):2493–2537.
- Jin Cong and Haitao Liu. 2014. Approaching human language with complex networks. *Physics of Life Reviews* 11(4):598 – 618. <https://doi.org/10.1016/j.plrev.2014.04.004>.
- Andre L. da Cunha, Lucilene B. de Sousa, Letícia L. Mansur, and Sandra M. Aluísio. 2015. Automatic proposition extraction from dependency trees: Helping early prediction of alzheimer’s disease from narratives. In *Proceedings of the 28th International Symposium on Computer-Based Medical Systems*. Institute of Electrical and Electronics Engineers, pages 127–130. <https://doi.org/10.1109/CBMS.2015.19>.
- Weissenbacher Davy, Johnson A. Travis, Wojtulewicz Laura, Dueck Amylou, Locke Dona, Caselli Richard, and Gonzalez Graciela. 2016. Towards automatic detection of abnormal cognitive decline and dementia through linguistic analysis of writing samples. In *Proceedings of the 15th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, pages 1198–1207. <https://doi.org/10.18653/v1/N16-1143>.
- Henrique F. de Arruda, Luciano F. Costa, and Diego R. Amancio. 2016. Using complex networks for text classification: Discriminating informative and imaginative documents. *EPL (Europhysics Letters)* 113(2):28007. <https://doi.org/10.1209/0295-5075/113/28007>.
- Martin Dyrba, Frederik Barkhof, Andreas Fellgiebel, Massimo Filippi, Lucrezia Hausner, Karlheinz Hauenstein, Thomas Kirste, and Stefan J. Teipel. 2015. Predicting prodromal alzheimer’s disease in subjects with mild cognitive impairment using machine learning classification of multimodal multicenter diffusion-tensor and magnetic resonance imaging data. *Journal of Neuroimaging* 25(5):738–747. <https://doi.org/10.1111/jon.12214>.
- Kathleen C. Fraser, Jed A. Meltzer, Naida L. Graham, Carol Leonard, Graeme Hirst, Sandra E. Black, and Elizabeth Rochon. 2014. Automated classification of primary progressive aphasia subtypes from narrative speech transcripts. *Cortex* 55:43–60. <https://doi.org/10.1016/j.cortex.2012.12.006>.
- Kathleen C. Fraser, Jed A. Meltzer, and Frank Rudzicz. 2015. Linguistic features identify alzheimer’s disease in narrative speech. *Journal of Alzheimer’s Disease* 49(2):407–422. <https://doi.org/10.3233/JAD-150520>.
- Peter Garrard, Vassiliki Rentoumi, Benno Gesierich, Bruce Miller, and Maria L. Gorno-Tempini. 2014. Machine learning approaches to diagnosis and laterality effects in semantic dementia discourse. *Cortex* 55:122–129. <https://doi.org/10.1016/j.cortex.2013.05.008>.
- Harold Goodglass, Edith Kaplan, and Barbara Barresi. 2001. *The Assessment of Aphasia and Related Disorders*. The Assessment of Aphasia and Related Disorders. Lippincott Williams & Wilkins.
- Arthur C. Graesser, Danielle S. McNamara, Max M. Louwerse, and Zhiqiang Cai. 2004.

- Coh-matrix: Analysis of text on cohesion and language. *Behavior research methods, instruments, & computers* 36(2):193–202. <https://doi.org/10.3758/BF03195564>.
- Ramon F. i Cancho, Ricard V. Solé, and Reinhard Köhler. 2004. Patterns in syntactic dependency networks. *Physical Review E* 69(5):051915. <https://doi.org/10.1103/PhysRevE.69.051915>.
- Ramon F. i Cancho and Richard V. Solé. 2001. The small world of human language. *Proceedings of the Royal Society of London B: Biological Sciences* 268(1482):2261–2265. <https://doi.org/10.1098/rspb.2001.1800>.
- William L. Jarrold, Bart Peintner, David Wilkins, Dimitra Vergryi, Colleen Richey, Maria L. Gorno-Tempini, and Jennifer Ogar. 2014. Aided diagnosis of dementia type through computer-based analysis of spontaneous speech. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics Workshop on Computational Linguistics and Clinical Psychology*. Association for Computational Linguistics, pages 27–36.
- William L. Jarrold, Bart Peintner, Eric Yeh, Ruth Krasnow, Harold S. Javitz, and Gary E. Swan. 2010. Language analytics for assessing brain health: Cognitive impairment, depression and pre-symptomatic alzheimer’s disease. In Yiyu Yao, Ron Sun, Tomaso Poggio, Jiming Liu, Ning Zhong, and Jimmy Huang, editors, *Proceedings of International Conference on Brain Informatics (BI 2010)*, Springer Berlin Heidelberg, pages 299–307. https://doi.org/10.1007/978-3-642-15314-3_28.
- Edith Kaplan, Harold Googlass, and Sandra Weintrab. 2001. *Boston naming test*. Lippincott Williams & Wilkins.
- A. Kertesz. 1982. *Western Aphasia Battery test manual*. Grune & Stratton.
- Maider Lehr, Emily T. Prud’hommeaux, Izhak Shafran, and Brian Roark. 2012. Fully automated neuropsychological assessment for detecting mild cognitive impairment. In *Proceedings of the 13th Annual Conference of the International Speech Communication Association*. pages 1039–1042.
- Jeaneth Machicao, Edilson A. Corrêa Jr, Gisele H. B. Miranda, Diego R. Amancio, and Odemir M. Bruno. 2016. Authorship attribution based on life-like network automata. *arXiv preprint arXiv:1610.06498*.
- Brian MacWhinney. 2000. *The CHILDES Project: Tools for analyzing talk*. Lawrence Erlbaum Associates, 3 edition.
- Rada Mihalcea and Dragomir Radev. 2011. *Graph-based natural language processing and information retrieval*. Cambridge University Press.
- Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. 2013a. Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 27th Annual Conference on Neural Information Processing Systems*. pages 3111–3119.
- Weerasak Muangpaisan, Chonachan Petcharat, and Varalak Srinonprasert. 2012. Prevalence of potentially reversible conditions in dementia and mild cognitive impairment in a geriatric clinic. *Geriatrics & gerontology international* 12(1):59–64. <https://doi.org/10.1111/j.1447-0594.2011.00728.x>.
- Sylvester O. Orimaye, Jojo Wong, and K. Jennifer Golden. 2014. Learning predictive linguistic features for alzheimer’s disease and related dementias using verbal utterances. In *Proceedings of the 1st Workshop on Computational Linguistics and Clinical Psychology (CLPsych)*. Association for Computational Linguistics, pages 78–87. www.aclweb.org/anthology/W/W14/W14-3210.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* 12:2825–2830.
- Bryan Perozzi, Rami Al-Rfou, Vivek Kulkarni, and Steven Skiena. 2014. Inducing language networks from continuous space word representations. In *Proceedings of the 5th Workshop on Complex Networks CompleNet 2014*, Springer, pages 261–273. https://doi.org/10.1007/978-3-319-05401-8_25.
- Ronald C. Petersen. 2004. Mild cognitive impairment as a diagnostic entity. *Journal of internal medicine* 256(3):183–194. <https://doi.org/10.1111/j.1365-2796.2004.01388.x>.
- Emily T. Prud’hommeaux and Brian Roark. 2011. Alignment of spoken narratives for automated neuropsychological assessment. In *Proceedings of Workshop on Automatic Speech Recognition & Understanding, ASRU*. Institute of Electrical and Electronics Engineers, pages 484–489. <https://doi.org/10.1109/ASRU.2011.6163979>.
- Brian Roark, Margaret Mitchell, John-Paul Hosom, Kristy Hollingshead, and Jeffrey Kaye. 2011. Spoken language derived measures for detecting mild cognitive impairment. *Transactions on Audio, Speech, and Language Processing, Institute of Electrical and Electronics Engineers* 19(7):2081–2090. <https://doi.org/10.1109/TASL.2011.2112351>.
- Ranzivelle M. Roxas and Giovanni Tapang. 2010. Prose and poetry classification and boundary detec-

tion using word adjacency network analysis. *International Journal of Modern Physics C* 21(04):503–512. <https://doi.org/10.1142/S0129183110015257>.

Eleanor M. Saffran, Rita S. Berndt, and Myrna F. Schwartz. 1989. The quantitative analysis of agrammatic production: Procedure and data. *Brain and language* 37(3):440–479. [https://doi.org/10.1016/0093-934X\(89\)90030-8](https://doi.org/10.1016/0093-934X(89)90030-8).

Thiago C. Silva and Diego R. Amancio. 2012. Word sense disambiguation via high order of learning in complex networks. *EPL (Europhysics Letters)* 98(5):58001.

Camila V. Teixeira, Lilian T. Gobbi, Danilla I. Corazza, Florindo Stella, José L. Costa, and Sebastião Gobbi. 2012. Non-pharmacological interventions on cognitive functions in older people with mild cognitive impairment (mci). *Archives of gerontology and geriatrics* 54(1):175–180. <https://doi.org/10.1016/j.archger.2011.02.014>.

László Tóth, Gábor Gosztolya, Veronika Vincze, Ildikó Hoffmann, and Gréta Szatlóczki. 2015. Automatic detection of mild cognitive impairment from spontaneous speech using asr. In *Proceedings of the 16th Annual Conference of the International Speech Communication Association*. International Speech and Communication Association, pages 2694–2698.

Marcos V. Treviso, Christopher Shulby, and Sandra M. Aluísio. 2017. Sentence segmentation in narrative transcripts from neuropsychological tests using recurrent convolutional neural networks. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, pages 1–10. <https://arxiv.org/abs/1610.00211>.

Veronika Vincze, Gábor Gosztolya, László Tóth, Ildikó Hoffmann, and Gréta Szatlóczki. 2016. Detecting mild cognitive impairment by exploiting linguistic information from transcripts. In *Proceedings of the 54th Annual Meeting of the Association Computer Linguistics*. Association for Computational Linguistics. <https://doi.org/10.18653/v1/P16-2030>.

Alyssa Weakley, Jennifer A. Williams, Maureen Schmitter-Edgecombe, and Diane J. Cook. 2015. Neuropsychological test selection for cognitive impairment classification: A machine learning approach. *Journal of clinical and experimental neuropsychology* 37(9):899–916. <https://doi.org/10.1080/13803395.2015.1067290>.

David Wechsler et al. 1997. *Wechsler memory scale (WMS-III)*. Psychological Corporation.

Zhi-Hua Zhou. 2012. *Ensemble methods: foundations and algorithms*. Chapman & Hall/CRC, 1st edition.

A Supplementary Material

Figure 3 is Cookie Theft picture, which was used in DementiaBank project.

Figure 4 is a sequence of pictures from the Cinderella story, which were used to elicit speech narratives.

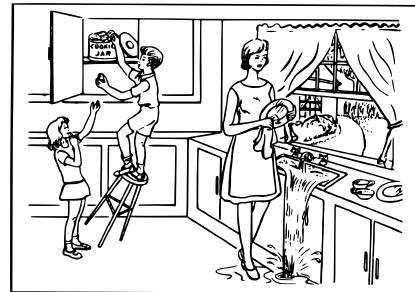


Figure 3: The Cookie Theft Picture, taken from the Boston Diagnostic Aphasia Examination (Goodglass et al., 2001).

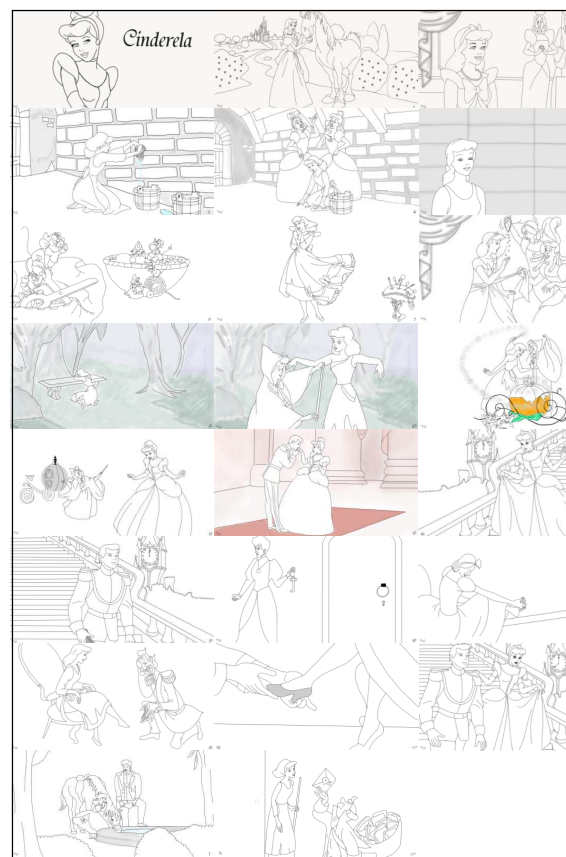


Figure 4: Sequence of Pictures of the of Cinderella story.

A.1 Examples of transcriptions

Below follows an example of a transcript of the Cookie Theft dataset.

You just want me to start talking ? Well the little girl is asking her brother we 'll say for a cookie . Now he 's getting the cookie one for him and one for her . He unbalances the step the little stool and he 's about to fall . And the lid 's off the cookie jar . And the mother is drying the dishes abstractly so she 's left the water running in the sink and it is spilling onto the floor . And there are two there 's look like two cups and a plate on the sink and board . And that boy 's wearing shorts and the little girl is in a short skirt . And the mother has an apron on . And she 's standing at the window . The window 's opened . It must be summer or spring . And the curtains are pulled back . And they have a nice walk around their house . And there 's this nice shrubbery it appears and grass . And there 's a big picture window in the background that has the drapes pulled off . There 's a not pulled off but pulled aside . And there 's a tree in the background . And the house with the kitchen has a lot of cupboard space under the sink board and under the cabinet from which the cookie you know cookies are being removed .

Below follows an excerpt of a transcript of the Cinderella dataset.

Original transcript in Portuguese:

ela morava com a madrasta as irmã né e ela era diferenciada das três era maltratada ela tinha que fazer limpeza na casa toda no castelo alias e as irmãs não faziam nada até que um dia chegou um convite do rei ele ia fazer um baile e a madrasta então é colocou que todas as filhas elas iam menos a cinderela bom como ela não tinha o vestido sapato as coisas tudo então ela mesmo teve que fazer a roupa dela começou a fazer ...

Translation of the transcript in English:

she lived with the stepmother the sister right and she was differentiated from the three was mistreated she had to do the cleaning in the entire house actually in the castle and the sisters didn't do anything until one day the king's invitation arrived he would invite everyone to a ball and then the stepmother is said that all the daughters they would go except for cinderella well since she didn't have a dress shoes all the things she had to make her own clothes she started to make them ...

A.2 Coh-Matrix-Dementia metrics

1. **Ambiguity:** verb ambiguity, noun ambiguity, adjective ambiguity, adverb ambiguity;
2. **Anaphoras:** adjacent anaphoric references,

anaphoric references;

3. **Basic Counts:** Flesch index, number of word, number of sentences, number of paragraphs, words per sentence, sentences per paragraph, syllables per content word, verb incidence, noun incidence, adjective incidence, adverb incidence, pronoun incidence, content word incidence, function word incidence;
4. **Connectives:** connectives incidence, additive positive connectives incidence, additive negative connectives incidence, temporal positive connectives incidence, temporal negative connectives incidence, casual positive connectives incidence, casual negative connectives incidence, logical positive connectives incidence, logical negative connectives incidence;
5. **Co-reference Measures:** adjacent argument overlap, argument overlap, adjacent stem overlap, stem overlap, adjacent content word overlap;
6. **Content Word Frequencies:** Content words frequency, minimum among content words frequency;
7. **Hypernyms:** Mean hypernyms per verb;
8. **Logic Operators:** Logic operators incidence, and incidence, or incidence, if incidence, negation incidence;
9. **Latent Semantic Analysis (LSA):** Average and standard deviation similarity between pairs of adjacent sentences in the text, Average and standard deviation similarity between all sentence pairs in the text, Average and standard deviation similarity between pairs of adjacent paragraphs in the text, Givenness average and standard deviation of each sentence in the text;
10. **Semantic Density:** content density;
11. **Syntactical Complexity:** only cross entropy;
12. **Tokens:** personal pronouns incidence, type-token ratio, Brunet index, Honoré Statistics.