# Language Networks: a Practical Approach

Jorge A. V. Tohalino[1] and Diego R. Amancio[1]

[1]*Institute of Mathematics and Computer Science,*

*University of São Paulo, São Carlos, Brazil*

(Dated: October 15, 2020)

## Abstract

This manuscript provides a short and practical introduction to the topic of language networks. This text aims at assisting researchers with no practical experience in text and/or network analysis. We provide a practical tutorial on how to model and characterize texts using network-based features. In this tutorial, we also include examples of pre-processing and network representations. A brief description of the main tasks allying network science and text analysis is also provided. A further development of this text shall include a practical description of network classification via machine learning methods.

# I. LANGUAGE NETWORKS: APPLICATIONS

Complex networks have been used to model a wide range of applications [1]. Particularly, this model have been used also to analyze and classify texts. From the theoretical perspective, complex networks have been used to analyze many aspects of human language, including the emergence of Zipf's Law [2] and semantic and cognitive properties related to language. Many of such analysis are in the Physics field. Practical approaches includes many applications in Computer Science. More recently, many applications of language network analysis. For reference purpose, we list below some of the applications of language network analysis. We include both theoretical and perspectives approaches. A more detailed account of network-based models can be found e.g. in [3].

1. *Authorship attribution and style attribution*: this task aims at identifying the authorship of a text, given a set of possible authors in the supervised case. Examples of network-based works addressing this topic include refs. [4–12].

2. *Stylometry*: this taks aims at characterizing different styles observed in literary movements, scientific journals and other issues. This topic is also closely related to the authorship attribution task, since authors have different writing style. Stylometric features can also be used e.g. to discriminate real from fake texts. Examples of network-based works addressing this topic include refs. [13–23].

3. *Complexity analysis*: this task aims at analyzing the features responsible for making a text complex. Research on this topic includes automatic essay scoring and comparison of language complexity patterns. Examples of network-based works addressing this topic includes [24–29].

4. *Document Summarization*: in this task one aims at identifying the most important information from large documents. In recent years complex networks have been used to identify the relevance of words, sentences or paragraphs via network centrality measurements. Examples of network-based works addressing this topic includes [30–35].

5. *Language theory*: this line of research includes the study of many aspects of linguistic theory, including the emergence of linguistic and complex systems patterns. Examples

of network-based works addressing this topic includes [24, 36–64].

6. *Semantic analysis*: complex networks have also been used to analyze the semantic aspects of texts. This is possible because different models can capture different language aspects [65–74].

7. *Story flow*: this line of studies how stories unfolds via network analysis. Examples of network-based works addressing this topic includes [75–78]. An example of story flow visualization is provided in Figure 1.

8. *Language and cognition*: many cognitive aspects can be studied via language. This includes studies the analysis the relationship between language features and cognitive impairment. Examples of studies on this topic can be found in refs. [48, 79–89].

9. *Keyword extraction*: this task is one of the most important and basic tasks in information retrieval. The identification of relevant words can also be important as a sub-task in other tasks, such as in document summarization and document clustering. Examples of studies on this topic can be found in refs. [34, 90–99].

10. *Sentiment analysis*: this task uses natural language processing and related textual approach to identify, in a systematic way, affective states in particular words, sentences or larger chunks of texts [100, 101].

## II. A PRACTICAL EXAMPLE: MODELING AND CHARACTERIZING CO-OCCURRENCE NETWORKS

A simple approach to model a text as a complex networks consists in linking adjacent words whenever they appear adjacent in the text. It has been shown that this type of representation captures mostly stylistic features [94]. While other types of representation are able of capturing semantic features (see Figure 1), our focus here is on co-occurrence networks. Word adjacency networks can be viewed as a simplification of text networks created via syntactical links, where nodes are words and edges are syntactical dependencies [102]. Because most of syntactical links are between adjacent words, a co-occurrence network can be viewed as a simplification of the syntactical model [102].
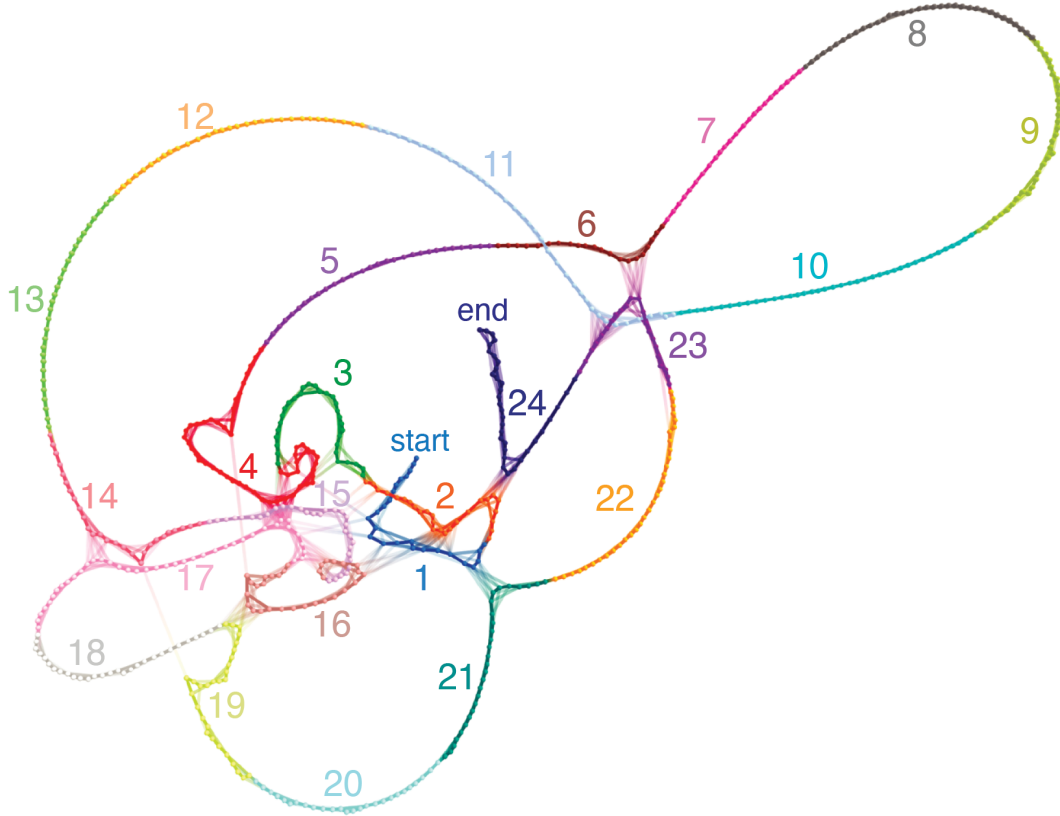
FIG. 1. Example of mesoscopic network, where nodes represent a set of consecutive paragraphs and edges are established between semantically similar nodes. The network was obtained from The Odyssey book series (attributed to Homer). Nodes colors and respective numbers indicate the books they belong to.

The basic steps in the creation of a word adjacency network include the removal of stopwords and lemmatization. The removal of stopwords in most applications because stopwords may largely affect the structure of networks. In most applications, because they mostly link content words, they can be replaced by edges. However, in some applications, stopwords may play a important role in detecting stylometric features [12]. After stopwords are removed, the remaining words are lemmatized so that verbs and nouns are mapped into their infinitive and singular forms, respectively.

In the first subsection of this document we describe how stopwords are removed. The lemmatization process is also detailed in the respective section. We then show how the pre-processed text is mapped into a co-occurrence network.

## A.  Pre-processing

Natural Language Toolkit (NLTK) is a well-known text processing Python tool [103]. In this example we use NLTK to import a list of English stopwords (function words). Note that this is a pre-selected list of words. Some strategies to define stopwords in any sequence of symbols can be also used for this purpose [94, 104]. As first step, we import libraries to handle stopwords and to perform the lemmatization process:

```python
import nltk
from nltk.corpus import stopwords
from nltk import word_tokenize
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet
```

FIG. 2. Libraries required to perform the removal of stopwords and the lemmatization process. To split the text into words we use the word tokenize library.

The first step in processing the text consists in tokenizing the text so that a list of words is provided as output:

```python
text = 'the children were playing games'
text_list = word_tokenize(text)
print('Original text:',text_list)
```

FIG. 3. The input text "the children were playing games" is tokenized to form a list of words using the function *word_tokenize*.

The output confirms that the text becomes a list of tokens: "the", "children", "were", "playing" and "games". In the next step, tokens identified as stopwords are removed from the text (see Figure 4). The output of Figure 4 lists the words "children", "playing" and "games", which confirms that "were" and "the" were identified as stopwords and removed. Note that the original word order is kept even after stopwords are disregarded. This is a essential feature in the construction and analysis of co-occurrence networks.

After stopwords are removed, *WordNetLemmatizer* in Figure 5 is used to extract the lemma of the remaining words. Here we do not provide the part-of-speech for each word, therefore all words are regarded as noun. The output of the code in Figure 5 lists the words "child", "playing" and "game".

```
def remove_stopwords(text_list):
    stop_list = stopwords.words('english')
    return [word for word in text_list if word not in stop_list]

text_list = remove_stopwords(text_list)
print('Text without stopwords:',text_list)
```

FIG. 4. Removing stopwords from the text.

```
lemmatized = [WordNetLemmatizer().lemmatize(word) for word in text_list]
print('Lemmatization using default tag (NOUN):', lemmatized)
```

FIG. 5. Example of word lemmatization using *WordNetLemmatizer*. The input parameters are the word and its corresponding part-of-speech. Because we did not provide the part-of-speech for each word, all words were regarded as noun.

The output reveals that the lemmatization process is not conducted with success for all words. Note that the word "playing" is not lemmatized since its part-of-speech is incorrect. In order to provide the correct part-of-speech for each word, we use *POSTagging* from NLTK, as shown in Figure 6. The output lists the part-of-speech of each word in a list format, i.e. ('children', 'NNS'), ('playing', 'VBG') and ('games', 'NNS').

```
tagged_words = nltk.pos_tag(text_list)
print('POSTags of words from text:',tagged_words)
```

FIG. 6. Example of part-of-speech tags obtained for the considered input words.

In order to use the obtained tags in *WordNetLemmatizer*, the following tags should be provided: "N", "V", "J" and "R" respectively for nouns, verbs, adjectives and adverbs. To map the output of the function used for pos-tagging into "N", "V", "J" and "R", we can use Figures 7 and 8. In particular, Figure 8 uses the code in Figure 7 to convert part-of-speech tags generated by NLTK to the WordNet format. The ouptut of Figure 8 provides the following list of tags: 'n', 'v' and 'n'. Finally, once the obtained tags are mapped to the desired format, they can be used to perform the lemmatization process using Figure 8. Now the output lists the words "child", "play" and "game" (in this order), as expected.

We summarize in Figure 9 all the steps required to pre-process a text before it can be used to create co-occurrence networks. The steps include (i) removal of stopwords (this is optional, depending on the model being used), (ii) part-of-speech tagging, and (iii) lemmatization.

6

```python
def get_wordnet_pos(tag):
    if tag.startswith('J'):
        return wordnet.ADJ
    elif tag.startswith('V'):
        return wordnet.VERB
    elif tag.startswith('N'):
        return wordnet.NOUN
    elif tag.startswith('R'):
        return wordnet.ADV
    else:
        return wordnet.NOUN


def convert_to_wn(tag):
    return get_wordnet_pos(tag)
```

FIG. 7. Functions used for generating "N", "V", "J" and "R" respectively for nouns, verbs, adjectives and adverbs.

```python
def convert_to_wn(tag):
    return get_wordnet_pos(tag)

wordnet_tags = [convert_to_wn(tag[1]) for tag in tagged_words]
print('Wordnet tags:', wordnet_tags)

lemmatized = [WordNetLemmatizer().lemmatize(word, tag) for  word,
tag in zip(text_list, wordnet_tags)]
print('Lemmatization using word tags:', lemmatized)
```

FIG. 8. Converting part-of-speech tags and performing the lemmatization process.

## B. Creating a co-occurrence text-network

Co-occurrence networks are well-known network models used to represent texts as complex networks. This model represents nodes as words and edges are established between adjacent (or nearby) words. Some works use windows so that words in the same window are linked as a clique. In the previous sentence, for example, the three first words would generate the following edges: "some" – "works", "works" – "use" and "some" – "use" (we disregarded the pre-processing steps in this example). Usually, the window size ($w$) ranges between 1 and 3.

Before implementing a function to create a network from a text, we create an auxiliary function *get_neighbors*. It returns the $w$ nearest neighbors of a word in a word list. This function takes into account both left and right neighbors. This implementation is provided

7

```
def text_processing(text):
    text_list = word_tokenize(text)
    text_list = remove_stopwords(text_list)
    tagged_words = nltk.pos_tag(text_list)
    wordnet_tags = [convert_to_wn(tag[1]) for tag in tagged_words]
    lemmatized = [WordNetLemmatizer().lemmatize(word, tag) for
    word, tag in zip(text_list, wordnet_tags)]
    return lemmatized
```

FIG. 9. Function used to pre-process a text. The pre-processing comprises the steps of (i) removing stopwords; (ii) part-of-speech tagging; and (iii) lemmatization. In some network models, the removal of stopwords is an optional step.

in Figure 10. *CNetwork* – the class responsible for mapping a text into a network – has as inputs the raw text and the window size $w$. The function returns the obtained network, which is represented using igraph library [105] (see Figure 11). *get_network()* creates a adjacency matrix in order to store adjacency information. Then the function seeks the $w$ nearest neighbors for each word. This is used to define the neighbor nodes for each word. The obtained matrix is used as input to the function *igraph.Graph.Adjacency()*, which returns the network representation in the igraph format.

```
import numpy as np
import igraph
from nltk import word_tokenize

def get_neighbors(word_list, index, w):
    if index - w >= 0:
        left = word_list[index - w:index]
    else:
        left = word_list[:index]
    right = word_list[index + 1:index + 1 + w]
    return set(left + right)
```

FIG. 10. This function take as input parameters a list of words (*word_list*), the index of the target word in the list (*index*), and the window size $w$.

The application of *CNetwork* is illustrated in Figure 12. We start by creating a new text that shall be mapped into a co-occurrence network. The tokenization process in Figure 12 provides the following list of words: 'today', 'we', 'are', 'learning', 'some', 'concepts', 'of', 'complex', 'networks', 'and', 'machine' and 'learning'. The tokenized text is then used as input to *cNetwork*. In this example, we use $w = 2$. The following result is obtained after modeling the text as a network:

```python
class CNetwork(object):

    def __init__(self, text, window):
        self.document = text
        self.words = set(self.document)
        self.vocab_index = {word:i for i, word in enumerate(self.words)}
        self.window = window

    def get_network(self):
        matrix = np.zeros((len(self.words), len(self.words)))
        for index, word in enumerate(self.document):
            neighbors = get_neighbors(self.document, index, self.window)
            word_index = self.vocab_index[word]
            for neighbor in neighbors:
                neighbor_index = self.vocab_index[neighbor]
                matrix[word_index][neighbor_index] = 1
        np.fill_diagonal(matrix, 0)
        network = igraph.Graph.Adjacency(matrix.tolist(), mode="undirected")
        network.vs['name'] = list(self.words)
        return network
```

FIG. 11. *CNetwork* creates a network model from a raw text. The function returns the obtained network,which is represented using igraph library.

1. Number of network nodes: 11.

2. Network nodes: 'of', 'networks', 'concepts', 'complex', 'and', 'machine', 'learning', 'some', 'are', 'we' and 'today'.

3. Number of edges: 21.

4. Edge list: (0, 1), (0, 2), (0, 3), (0, 7), (1, 3), (1, 4), (1, 5), (2, 3), (2, 6), (2, 7), (3, 4), (4, 5), (4, 6), (5, 6), (6, 7), (6, 8), (6, 9), (7, 8), (8, 9), (8, 10) and (9, 10).

In Figure 12, we also list some network measurements extracted from the obtained co-occurrence network. A detailed description of network measurements can be found in ref. [106]. The obtained values for the considered network measurements are:

1. Degree: 4, 4, 4, 3, 4, 3, 6, 4, 4, 2, 4.

2. Degree of some nodes: 4, 4, 6.

3. PageRank: 0.09, 0.1, 0.09, 0.07, 0.09, 0.08, 0.14, 0.09, 0.09, 0.06, 0.09.

4. Betweenness: 3.82, 6.27, 1.5, 1.73, 2.0, 2.57, 17.47, 5.37, 2.95, 0.0, 3.33.

```python
text = 'today we are learning some concepts
of complex networks and machine learning '
text = word_tokenize(text)
print('Source text:', text)

obj = CNetwork(text, 2)
network = obj.get_network()
print('Number of network nodes:',network.vcount())
print('Network nodes:', network.vs['name'])
print('Number of edges:', len(network.get_edgelist()))
print('Edge list:',network.get_edgelist())

degree = network.degree()
degree_some_nodes = network.degree(['complex', 'networks', 'learning'])
pageRank = network.pagerank()
pageRank = [round(pr,2) for pr in pageRank]
betweenness = network.betweenness()
betweenness = [round(btw,2) for btw in betweenness]
print('Degree:', degree)
print('Degree of some nodes:', degree_some_nodes)
print('PageRank:', pageRank)
print('Betweenness:', betweenness)
```

FIG. 12. Example of network measurements extracted from the co-occurrence text network.

## III.   FINAL REMARKS

This text is intended to be a didactic approach for those interested in working at the intersection of network science, natural language processing, pattern recognition and applications. A further development of this manuscript will include a practical description of network characterization and classification. This further development shall touch, therefore, many of the applications described in Section I. We also intend to include additional models of network representation using advanced recent neural networks embeddings (see e.g. refs. [67, 107]).

[1] F. Battiston, G. Cencetti, I. Iacopini, V. Latora, M. Lucas, A. Patania, J.-G. Young, and G. Petri, "Networks beyond pairwise interactions: structure and dynamics," *Physics Reports*, 2020.

[2] R. Ferrer-i Cancho, C. Bentz, and C. Seguin, "Optimal coding and the origins of zipfian laws," *Journal of Quantitative Linguistics*, pp. 1–30, 2020.

[3] J. Cong and H. Liu, "Approaching human language with complex networks," *Physics of life reviews*, vol. 11, no. 4, pp. 598–618, 2014.

[4] S. Segarra, M. Eisen, and A. Ribeiro, "Authorship attribution through function word adjacency networks," *IEEE Transactions on Signal Processing*, vol. 63, no. 20, pp. 5464–5478, 2015.

[5] H. F. Arruda, L. F. Costa, and D. R. Amancio, "Using complex networks for text classification: Discriminating informative and imaginative documents," *EPL (Europhysics Letters)*, vol. 113, no. 2, p. 28007, 2016.

[6] V. Q. Marinho, G. Hirst, and D. R. Amancio, "Authorship attribution via network motifs identification," in *2016 5th Brazilian Conference on Intelligent Systems (BRACIS)*, pp. 355–360, IEEE, 2016.

[7] C. Akimushkin, D. R. Amancio, and O. N. Oliveira Jr, "On the role of words in the network structure of texts: Application to authorship attribution," *Physica A: Statistical Mechanics and its Applications*, vol. 495, pp. 49–58, 2018.

[8] J. Machicao, E. A. Corrêa Jr, G. H. Miranda, D. R. Amancio, and O. M. Bruno, "Authorship attribution based on life-like network automata," *PloS one*, vol. 13, no. 3, p. e0193703, 2018.

[9] C. Akimushkin, D. R. Amancio, and O. N. Oliveira Jr, "Text authorship identified using the dynamics of word co-occurrence networks," *PloS one*, vol. 12, no. 1, p. e0170527, 2017.

[10] D. R. Amancio, "Authorship recognition via fluctuation analysis of network topology and word intermittency," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2015, no. 3, p. P03005, 2015.

[11] S. Segarra, M. Eisen, G. Egan, and A. Ribeiro, "Attributing the authorship of the henry vi plays by word adjacency," *Shakespeare Quarterly*, vol. 67, no. 2, pp. 232–256, 2016.

[12] S. Segarra, M. Eisen, and A. Ribeiro, "Authorship attribution using function words adjacency networks," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 5563–5567, IEEE, 2013.

[13] M. Eisen, A. Ribeiro, S. Segarra, and G. Egan, "Stylometric analysis of early modern period english plays," *Digital Scholarship in the Humanities*, vol. 33, no. 3, pp. 500–528, 2018.

[14] V. Q. Marinho, G. Hirst, and D. R. Amancio, "Labelled network subgraphs reveal stylistic subtleties in written texts," *Journal of Complex Networks*, vol. 6, no. 4, pp. 620–638, 2018.

[15] D. R. Amancio, "A complex network approach to stylometry," *PloS one*, vol. 10, no. 8, p. e0136076, 2015.

[16] D. R. Amancio, "Probing the topological properties of complex networks modeling short written texts," *PloS one*, vol. 10, no. 2, p. e0118394, 2015.

[17] D. R. Amancio, O. N. Oliveira Jr, and L. d. F. Costa, "Structure–semantics interplay in complex networks and its effects on the predictability of similarity in texts," *Physica A: Statistical Mechanics and its Applications*, vol. 391, no. 18, pp. 4406–4419, 2012.

[18] D. R. Amancio, "Comparing the topological properties of real and artificially generated scientific manuscripts," *Scientometrics*, vol. 105, no. 3, pp. 1763–1779, 2015.

[19] J. Jiang, W. Yu, and H. Liu, "Does scale-free syntactic network emerge in second language learning?," *Frontiers in Psychology*, vol. 10, p. 925, 2019.

[20] T. Stanisz, J. Kwapień, and S. Drożdż, "Linguistic data mining with complex networks: a stylometric-oriented approach," *Information Sciences*, vol. 482, pp. 301–320, 2019.

[21] R. M. Roxas-Villanueva, M. K. Nambatac, and G. Tapang, "Characterizing english poetic style using complex networks," *International Journal of Modern Physics C*, vol. 23, no. 02, p. 1250009, 2012.

[22] J. Stevanak, D. M. Larue, and L. D. Carr, "Distinguishing fact from fiction: Pattern recognition in texts using complex networks," *arXiv preprint arXiv:1007.3254*, 2010.

[23] R. M. Roxas and G. Tapang, "Prose and poetry classification and boundary detection using word adjacency network analysis," *International Journal of Modern Physics C*, vol. 21, no. 04, pp. 503–512, 2010.

[24] H. Liu and W. Li, "Language clusters based on linguistic complex networks," *Chinese Science Bulletin*, vol. 55, no. 30, pp. 3458–3465, 2010.

[25] H. Liu, "The complexity of chinese syntactic dependency networks," *Physica A: Statistical Mechanics and its Applications*, vol. 387, no. 12, pp. 3048–3058, 2008.

[26] D. R. Amancio, S. M. Aluisio, O. N. Oliveira Jr, and L. d. F. Costa, "Complex networks analysis of language complexity," *EPL (Europhysics Letters)*, vol. 100, no. 5, p. 58002, 2012.

[27] H. Liu and C. Xu, "Can syntactic networks indicate morphological complexity of a language?," *EPL (Europhysics Letters)*, vol. 93, no. 2, p. 28005, 2011.

[28] L. Antiqueira, M. d. G. V. Nunes, O. Oliveira Jr, and L. d. F. Costa, "Strong correlations between text quality and complex networks features," *Physica A: Statistical Mechanics and its Applications*, vol. 373, pp. 811–820, 2007.

[29] J. P. Cárdenas, I. González, G. Vidal, and M. A. Fuentes, "Does network complexity help organize babel's library?," *Physica A: Statistical Mechanics and its Applications*, vol. 447, pp. 188–198, 2016.

[30] J. V. Tohalino and D. R. Amancio, "Extractive multi-document summarization using multi-layer networks," *Physica A: Statistical Mechanics and its Applications*, vol. 503, pp. 526–539, 2018.

[31] J. V. Tohalino and D. R. Amancio, "Extractive multi-document summarization using dynamical measurements of complex networks," in *2017 Brazilian Conference on Intelligent Systems (BRACIS)*, pp. 366–371, IEEE, 2017.

[32] L. Antiqueira, O. N. Oliveira Jr, L. da Fontoura Costa, and M. d. G. V. Nunes, "A complex network approach to text summarization," *Information Sciences*, vol. 179, no. 5, pp. 584–599, 2009.

[33] R. Mihalcea, "Graph-based ranking algorithms for sentence extraction, applied to text summarization," in *Proceedings of the ACL interactive poster and demonstration sessions*, pp. 170–173, 2004.

[34] D. R. Amancio, M. G. Nunes, O. N. Oliveira Jr, and L. d. F. Costa, "Extractive summarization using complex networks and syntactic dependency," *Physica A: Statistical Mechanics and its Applications*, vol. 391, no. 4, pp. 1855–1864, 2012.

[35] A. Mehler, "Structural similarities of complex networks: A computational model by example of wiki graphs," *Applied Artificial Intelligence*, vol. 22, no. 7-8, pp. 619–683, 2008.

[36] M. Stella and M. Brede, "Patterns in the english language: phonological networks, percolation and assembly models," *Journal of Statistical Mechanics: Theory and Experiment*,

vol. 2015, no. 5, p. P05006, 2015.

[37] M. Stella, S. De Nigris, A. Aloric, and C. S. Siew, "Forma mentis networks quantify crucial differences in stem perception between students and experts," *PloS one*, vol. 14, no. 10, p. e0222870, 2019.

[38] M. Stella, "Modelling early word acquisition through multiplex lexical networks and machine learning," *Big Data and Cognitive Computing*, vol. 3, no. 1, p. 10, 2019.

[39] V. Bochkarev, A. Shevlyakova, and E. Y. Lerner, "Modelling of growth of syntactic relations network in english and russian," *SCOPUS17426588-2018-1141-1-SID85059371763*, 2018.

[40] R. F. I. Cancho and R. V. Solé, "The small world of human language," *Proceedings of the Royal Society of London. Series B: Biological Sciences*, vol. 268, no. 1482, pp. 2261–2265, 2001.

[41] R. F. Cancho and R. V. Solé, "Least effort and the origins of scaling in human language," *Proceedings of the National Academy of Sciences*, vol. 100, no. 3, pp. 788–791, 2003.

[42] C. T. Kello, G. D. Brown, R. F. Cancho, J. G. Holden, K. Linkenkaer-Hansen, T. Rhodes, and G. C. Van Orden, "Scaling laws in cognitive sciences," *Trends in cognitive sciences*, vol. 14, no. 5, pp. 223–232, 2010.

[43] R. F. Cancho, "Euclidean distance between syntactically linked words," *Physical Review E*, vol. 70, no. 5, p. 056135, 2004.

[44] R. Ferrer Cancho, "Why do syntactic links not cross?," *Europhysics Letters*, vol. 76, no. 6, pp. 1228–1235, 2006.

[45] R. F. Cancho, A. Capocci, and G. Caldarelli, "Spectral methods cluster words of the same class in a syntactic dependency network," *International Journal of Bifurcation and Chaos*, vol. 17, no. 07, pp. 2453–2463, 2007.

[46] R. F. i Cancho, "The structure of syntactic dependency networks: insights from recent advances in network theory," *Problems of quantitative linguistics*, pp. 60–75, 2005.

[47] M. Choudhury and A. Mukherjee, "The structure and dynamics of linguistic networks," in *Dynamics on and of Complex Networks*, pp. 145–166, Springer, 2009.

[48] R. V. Solé, B. Corominas-Murtra, S. Valverde, and L. Steels, "Language networks: Their structure, function, and evolution," *Complexity*, vol. 15, no. 6, pp. 20–26, 2010.

[49] M. M. Soares, G. Corso, and L. Lucena, "The network of syllables in portuguese," *Physica A: Statistical Mechanics and its Applications*, vol. 355, no. 2-4, pp. 678–684, 2005.

[50] G. Peng, J. W. Minett, and W. S.-Y. Wang, "The networks of syllables and characters in chinese," *Journal of Quantitative Linguistics*, vol. 15, no. 3, pp. 243–255, 2008.

[51] J. Li and J. Zhou, "Chinese character structure analysis based on complex networks," *Physica A: Statistical Mechanics and its Applications*, vol. 380, pp. 629–638, 2007.

[52] Y. Li, L. Wei, Y. Niu, and J. Yin, "Structural organization and scale-free properties in chinese phrase networks," *Chinese Science Bulletin*, vol. 50, no. 13, pp. 1305–1309, 2005.

[53] S. Zhou, G. Hu, Z. Zhang, and J. Guan, "An empirical study of chinese language networks," *Physica A: Statistical Mechanics and its Applications*, vol. 387, no. 12, pp. 3039–3047, 2008.

[54] M. Brede and D. Newth, "Patterns in syntactic dependency networks from authored and randomised texts," *Complexity International*, vol. 12, no. msid23, 2008.

[55] L. Sheng and C. Li, "English and chinese languages as weighted complex networks," *Physica A: Statistical Mechanics and its Applications*, vol. 388, no. 12, pp. 2561–2570, 2009.

[56] I. Grabska-Gradzińska, A. Kulig, J. Kwapień, and S. Drożdż, "Complex network analysis of literary and scientific texts," *International Journal of Modern Physics C*, vol. 23, no. 07, p. 1250051, 2012.

[57] W. Liang, Y. Shi, K. T. Chi, and Y. Wang, "Study on co-occurrence character networks from chinese essays in different periods," *Science China Information Sciences*, vol. 55, no. 11, pp. 2417–2427, 2012.

[58] Y. Gao, W. Liang, Y. Shi, and Q. Huang, "Comparison of directed and weighted co-occurrence networks of six languages," *Physica A: Statistical Mechanics and its Applications*, vol. 393, pp. 579–589, 2014.

[59] R. Čech and J. Mačutek, "Word form and lemma syntactic dependency networks in czech: A comparative study," *Glottometrics*, vol. 19, pp. 85–98, 2009.

[60] C. X. X. C. L. Wenwen, "Extracting valency patterns of word classes from syntactic complex networks," *depling. org*, p. 165.

[61] O. Abramov and A. Mehler, "Automatic language classification by means of syntactic dependency networks," *Journal of Quantitative Linguistics*, vol. 18, no. 4, pp. 291–336, 2011.

[62] J. Ke and Y. Yao, "Analysing language development from a network approach," *Journal of Quantitative Linguistics*, vol. 15, no. 1, pp. 70–99, 2008.

[63] E. Gibson, "Linguistic complexity: Locality of syntactic dependencies," *Cognition*, vol. 68, no. 1, pp. 1–76, 1998.

[64] H. Liu, "Dependency direction as a means of word-order typology: A method based on dependency treebanks," *Lingua*, vol. 120, no. 6, pp. 1567–1578, 2010.

[65] D. R. Amancio, O. N. Oliveira Jr, and L. da F Costa, "Using complex networks to quantify consistency in the use of words," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2012, no. 01, p. P01004, 2012.

[66] E. A. Corrêa Jr and D. R. Amancio, "Word sense induction using word embeddings and community detection in complex networks," *Physica A: Statistical Mechanics and its Applications*, vol. 523, pp. 180–190, 2019.

[67] E. A. Corrêa Jr, V. Q. Marinho, and D. R. Amancio, "Semantic flow in language networks discriminates texts by genre and publication date," *Physica A: Statistical Mechanics and its Applications*, vol. 557, p. 124895, 2020.

[68] E. A. Correa Jr, A. A. Lopes, and D. R. Amancio, "Word sense disambiguation: A complex network approach," *Information Sciences*, vol. 442, pp. 103–113, 2018.

[69] H. Liu, "Statistical properties of chinese semantic networks," *Chinese Science Bulletin*, vol. 54, no. 16, pp. 2781–2785, 2009.

[70] R. F. Cancho and R. V. Solé, "Two regimes in the frequency of words and the origins of complex lexicons: Zipf's law revisited," *Journal of Quantitative Linguistics*, vol. 8, no. 3, pp. 165–173, 2001.

[71] J. Borge-Holthoefer and A. Arenas, "Semantic networks: Structure and dynamics," *Entropy*, vol. 12, no. 5, pp. 1264–1302, 2010.

[72] R. Mihalcea and D. Radev, *Graph-based natural language processing and information retrieval*. Cambridge university press, 2011.

[73] R. V. Solé and L. F. Seoane, "Ambiguity in language networks," *The Linguistic Review*, vol. 32, no. 1, pp. 5–35, 2015.

[74] T. C. Silva and D. R. Amancio, "Discriminating word senses with tourist walks in complex networks," *The European Physical Journal B*, vol. 86, no. 7, p. 297, 2013.

[75] D. R. Amancio, "Network analysis of named entity co-occurrences in written texts," *EPL (Europhysics Letters)*, vol. 114, no. 5, p. 58005, 2016.

[76] N. Dekker, T. Kuhn, and M. van Erp, "Evaluating social network extraction for classic and modern fiction literature," *PeerJ Preprints*, vol. 6, p. e27263v1, 2018.

[77] A. J. Holanda, M. Matias, S. M. Ferreira, G. M. Benevides, and O. Kinouchi, "Character networks and book genre classification," *International Journal of Modern Physics C*, vol. 30, no. 08, p. 1950058, 2019.

[78] S. D. Prado, S. R. Dahmen, A. L. Bazzan, P. M. Carron, and R. Kenna, "Temporal network analysis of literary texts," *Advances in Complex Systems*, vol. 19, no. 03, p. 1650005, 2016.

[79] L. B. dos Santos, E. A. C. Júnior, O. N. Oliveira Jr, D. R. Amancio, L. L. Mansur, and S. M. Aluísio, "Enriching complex networks with word embeddings for detecting mild cognitive impairment from speech transcripts.," in *ACL (1)*, pp. 1284–1296, 2017.

[80] M. Stella and Y. N. Kenett, "Viability in multiplex lexical networks and machine learning characterizes human creativity," *Big Data and Cognitive Computing*, vol. 3, no. 3, p. 45, 2019.

[81] M. Stella, "Cohort and rhyme priming emerge from the multiplex network structure of the mental lexicon," *Complexity*, vol. 2018, 2018.

[82] M. Stella and A. Zaytseva, "Forma mentis networks map how nursing and engineering students enhance their mindsets about innovation and health during professional growth," *PeerJ Computer Science*, vol. 6, p. e255, 2020.

[83] M. Choudhury, M. Thomas, A. Mukherjee, A. Basu, and N. Ganguly, "How difficult is it to develop a perfect spell-checker? a cross-linguistic analysis through complex network approach," *arXiv preprint physics/0703198*, 2007.

[84] A. Baronchelli, R. Ferrer-i Cancho, R. Pastor-Satorras, N. Chater, and M. H. Christiansen, "Networks in cognitive science," *Trends in cognitive sciences*, vol. 17, no. 7, pp. 348–360, 2013.

[85] M. Sigman and G. A. Cecchi, "Global organization of the wordnet lexicon," *Proceedings of the National Academy of Sciences*, vol. 99, no. 3, pp. 1742–1747, 2002.

[86] A. de Jesus Holanda, I. T. Pisa, O. Kinouchi, A. S. Martinez, and E. E. S. Ruiz, "Thesaurus as a complex network," *Physica A: Statistical Mechanics and its Applications*, vol. 344, no. 3-4, pp. 530–536, 2004.

[87] M. Steyvers and J. B. Tenenbaum, "The large-scale structure of semantic networks: Statistical analyses and a model of semantic growth," *Cognitive science*, vol. 29, no. 1, pp. 41–78, 2005.

[88] P. Gravino, V. D. Servedio, A. Barrat, and V. Loreto, "Complex structures and semantics in free word association," *Advances in Complex Systems*, vol. 15, no. 03n04, p. 1250054, 2012.

[89] M. Stella, "Text-mining forma mentis networks reconstruct public perception of the stem gender gap in social media," *PeerJ Computer Science*, vol. 6, p. e295, Sept. 2020.

[90] S. Lahiri, S. R. Choudhury, and C. Caragea, "Keyword and keyphrase extraction using centrality measures on collocation networks (2014)," *arXiv preprint arXiv:1401.6571*.

[91] Z. J. Zhan, F. Lin, and X. P. Yang, "Keyword extraction of document based on weighted complex network," in *Advanced Materials Research*, vol. 403, pp. 2146–2151, Trans Tech Publ, 2012.

[92] S. Beliga, A. Meštrović, and S. Martinčić-Ipšić, "Selectivity-based keyword extraction method," *International Journal on Semantic Web and Information Systems (IJSWIS)*, vol. 12, no. 3, pp. 1–26, 2016.

[93] R. Wang, W. Liu, and C. McDonald, "Using word embeddings to enhance keyword identification for scientific publications," in *Australasian Database Conference*, pp. 257–268, Springer, 2015.

[94] D. R. Amancio, E. G. Altmann, D. Rybski, O. N. Oliveira Jr, and L. d. F. Costa, "Probing the statistical properties of unknown texts: application to the voynich manuscript," *PLoS One*, vol. 8, no. 7, p. e67310, 2013.

[95] H. Li, H. An, Y. Wang, J. Huang, and X. Gao, "Evolutionary features of academic articles co-keyword network and keywords co-occurrence network: Based on two-mode affiliation network," *Physica A: Statistical Mechanics and its Applications*, vol. 450, pp. 657–669, 2016.

[96] S. Sišovic, S. Martincic-Ipšic, and A. Meštrovic, "Toward network-based keyword extraction from multitopic web documents," in *Proceedings of 6th International Conference on Information Technologies and Information Society (ITIS2014),'marješke toplice*, pp. 18–27.

[97] R. Yan and G.-l. Gao, "Chinese keywords identification based on strength entropy," *Computer Engineering & Science*, no. 11, p. 32, 2016.

[98] M. Garg and M. Kumar, "Identifying influential segments from word co-occurrence networks using ahp," *Cognitive Systems Research*, vol. 47, pp. 28–41, 2018.

[99] J. Mathiesen, L. Angheluta, and M. H. Jensen, "Statistics of co-occurring keywords on twitter," *arXiv preprint arXiv:1401.4140*, 2014.

[100] M. Mitrović, G. Paltoglou, and B. Tadić, "Networks and emotion-driven user communities at popular blogs," *The European Physical Journal B*, vol. 77, no. 4, pp. 597–609, 2010.

[101] P. Fornacciari, M. Mordonini, and M. Tomaiuolo, "Social network and sentiment analysis on twitter: Towards a combined approach.," in *KDWeb*, pp. 53–64, 2015.

[102] R. F. i Cancho, R. V. Solé, and R. Köhler, "Patterns in syntactic dependency networks," *Physical Review E*, vol. 69, no. 5, p. 051915, 2004.

[103] E. Loper and S. Bird, "Nltk: the natural language toolkit," *arXiv preprint cs/0205028*, 2002.

[104] C. Carretero-Campos, P. Bernaola-Galván, A. Coronado, and P. Carpena, "Improving statistical keyword detection in short texts: Entropic and clustering approaches," *Physica A: Statistical Mechanics and its Applications*, vol. 392, no. 6, pp. 1481–1492, 2013.

[105] G. Csardi, T. Nepusz, *et al.*, "The igraph software package for complex network research," *InterJournal, complex systems*, vol. 1695, no. 5, pp. 1–9, 2006.

[106] L. F. Costa, F. A. Rodrigues, G. Travieso, and P. R. Villas Boas, "Characterization of complex networks: A survey of measurements," *Advances in physics*, vol. 56, no. 1, pp. 167–242, 2007.

[107] L. V. Quispe, J. A. Tohalino, and D. R. Amancio, "Using virtual edges to improve the discriminability of co-occurrence text networks," *Physica A: Statistical Mechanics and its Applications*, p. 125344, 2020.