

# Lexical diversity and Mild Cognitive Impairment

Sofie Johansson<sup>1</sup>, Kristina Lundholm Fors<sup>1</sup>, Malin Antonsson<sup>1</sup>,  
Dimitrios Kokkinakis<sup>1,2</sup>

<sup>1</sup>Department of Swedish, University of Gothenburg, Sweden

<sup>2</sup>Centre for Ageing and Health, University of Gothenburg, Sweden

<https://doi.org/10.36505/ExLing-2019/10/0029/000391>

## Abstract

This paper explores the role that various lexical-based measures play for differentiating between individuals with mild forms of cognitive impairment (MCI) and healthy controls (HC). Recent research underscores the importance of language and linguistic analysis as essential components that can contribute to a variety of sensitive cognitive measures for the identification of milder forms of cognitive impairment. Subtle language changes serve as a sign that an individual's cognitive functions have been impacted, potentially leading to early diagnosis. Our research aims to identify linguistic biomarkers that could distinguish between individuals with MCI and HC and also be useful in predicting MCI.

Key words: mild cognitive impairment/MCI; lexical diversity; language; Swedish

## Introduction

The number of people living with dementia worldwide is projected to be 65.7 million in 2030, and this number will double every 20 years *cf.* Prince et al. (2013). It has been established that early diagnosis is beneficial and ultimately current research is exploring the possibility of identifying persons with mild forms of cognitive impairment at an early stage.

Mild cognitive impairment (MCI) is a condition characterized by cognitive decline greater than expected for an individual's age and education level. As the MCI progresses, MCI individuals face a higher risk of developing Alzheimer's Disease (AD). While language impairments have been well described in AD, language impairment in people with MCI is less well understood and the need for further research in all aspects of language and during all stages of the disease, has been recently emphasized in various studies; *cf.* Laske et al. 2015, Boschi et al. 2017; Beltrami et al. 2018.

In this paper, we examine lexical aspects of lexical diversity of individuals with MCI and age-matched controls and apply various measures described in relevant literature trying to identify whether there are correlations between lexical diversity measures and the participants' characteristics. The objective of the current study is to determine which and to what extent various lexical diversity measures can differentiate the two groups and what is the range of scores of the measures found in transcriptions of the oral narratives.

## Background

Multiple components of language can be assessed in order to reveal linguistic features that are likely to serve as discriminators between individuals with MCI and cognitively healthy controls. One of these components is at the lexical level, but research so far has shown conflicting results, probably because impairment may only appear at the *semantic* and *macrolinguistic* language levels on e.g. discourse relations (Masrani et al., 2017). Aramaki et al. (2016) used transcriptions of both written and spoken samples in order to measure various language ability scores. Their analysis of the spoken narrative showed e.g. that MCIs had a significantly larger vocabulary size which might indicate compensatory behavior for MCI persons. Fergadiotis et al. (2013) collected validity evidence regarding techniques for measuring lexical diversity for the study of aphasic discourse. Two of the tested lexical diversity scores, the Measure of Textual Lexical Diversity, MTLTD (McCarthy, Jarvis, 2010) and the Moving-Average Type-Token Ratio, MATTR (Covington, McFall, 2010) yielded the strongest evidence for producing unbiased lexical diversity scores, suggesting that they may be the best measures for lexical diversity in people with aphasia. For other relevant research on lexical profiles and linguistic features cf. Laufer & Nation (1995) and Biber (1995).

## Data, methods, results and limitations

Participants for this study were recruited from the longitudinal Gothenburg MCI study (Wallin et al., 2016). All subjects were native speakers of Swedish, and all studies are approved by the local ethical committee review board. In the current analyses, we include only participants who had completed the Cookie Theft task (Goodglass et al., 1983), a widely used test to elicit narrative speech. Participants were asked to describe everything they saw in the picture, and to talk for as long as they liked. The narratives were digitally recorded and manually transcribed according to a detailed protocol developed by the authors. Only the orthographic transcriptions of the spoken samples for each participant were used in the measurement of lexical diversity in this study. For each participant, we applied 16 lexical measures, e.g. long words, hapax legomena, frequency-related coverage at a general and a more specific level, originality, erroneous or non-existent words, contextual and non-contextual descriptors. As indicated by e.g. Laufer, Nation (1995), Richards et al. (2009) and Johansson, Ohlsson (2019), there are some lexical features which it is possible to relate to prominent or more proficient oral or written language, e.g. long words, hapax legomena and originality. The specific lexical measures in this work were selected since a larger vocabulary and richer language is commonly indicated by a higher lexical diversity, and contextual, rather than general language and vocabulary size, cf. Nation (2013) and Milton (2009). We assumed that there would be individual variables co-occurring as well as variables which would be significant to a particular group of individuals in the study. In order to

find possible correlations between the lexical measures, a bivariate 2-tailed correlation using Pearson's coefficient was conducted using IBM SPSS v 25.

The length of speech varied among the participants; the transcribed spoken picture descriptions ranged from 57 to 617 tokens (mean 185.47) for the HCs; and 46 to 481 tokens (mean 185.45) for the MCIs. The correlations which could be related to individual differences or similarities were found to be age-related to common words (.395\*\*) and the 1000 most common tokens used in a Swedish everyday context (.435\*\*) and negative correlation to long words and contextual descriptors. Other individual correlations were contextual descriptors which were related to hapax legomena (.439\*\*) and long words (.439\*\*). Weak correlations related to the two groups (HC, MCI), were contextual descriptors (.267\*) regarding the HC group, and erroneous words (.259\*) regarding the MCIs. This study has limitations. First, our study comprised relatively few participants, 29 HCs and 26 MCI. As this was not a very large sample, future studies with larger sample sizes are needed to verify our findings. Next, the lexical analysis carried out to examine differences in elicited verbal production originates from research in language development and language proficiency which might affect the generalizability of the results, although these measures have been applied to communication disorders such as aphasia. Preliminary findings indicate that there is a difference in how the individuals express themselves which might be related to age. It seems that older participants use more common words as opposed to more frequent use of contextual descriptors. This difference seems to be unrelated to numbers of education years. Further studies would involve a more sophisticated categorization of linguistically descriptive features, such as active-passive use of verbs, and also, lexical density where the number of content words is compared to functions words.

### Conclusions and future work

Our preliminary results show some differences between the scores and groups. The validity needs to be further investigated using more lexically-based measures of lexical diversity and richness in order to capture the degree of difference between the groups, by applying e.g. clustering. Further investigations are also required to shed light on the relationship between the lexical diversity measures and neuropsychological tests, and also apply them to lemmatized transcriptions. In the near future, we also plan to augment these measures e.g. MTLT, MATTR, and also compare the measures at 2 points in time, since exactly the same cohort has recently repeated the same task and have also undergone renewed neuropsychological assessments.

### Acknowledgements

This work has received support from *Riksbankens Jubileumsfond* – the Swedish Foundation for Humanities & Social Sciences; grant NHS 14-1761:1.

## References

- Aramaki, E., Shikata, S., Miyabe, M., Kinoshita, A. 2016. Vocabulary size in speech may be an early indicator of cognitive impairment. *PLOS ONE* 11(5).
- Beltrami, D., Gagliardi, G., Rossini Favretti, R., Ghidoni, E., Tamburini, F., Calzà, L. 2018. Speech analysis by NLP techniques: A possible tool for very early detection of cognitive decline? *Frontiers in Aging Neuroscience*. 10.
- Biber, D. 1995. *Dimensions of Register Variation: A Cross-Ling Comparison*. CUP.
- Boschi, V., Catricalà E., Consonni M., Chesi C., Moro A., Cappa S.F. 2017. Connected speech in neurodegenerative language disorders: a review. *Frontiers in Psychology*, 8:269.
- Covington, M., McFall, J. 2010. Cutting the gordian knot: The moving-average type token ratio (MATTR). *J. of Quan. Ling.* 17, 94-100.
- Fergadiotis, G., Wright, H.H., Westa, T.M. 2013. Measuring lexical diversity in narrative discourse of people with aphasia. *Am J Speech Lang Pathol.* 22, 2.
- Goodglass, H., Barresi, B., Kaplan, E. 1983. *Boston Diagn Aphasia Exam*. Kluwer.
- Johansson, S., Ohlsson, E. 2019. Visualizing vocabulary. *Multiling Matters*.
- Laske, C. et al. 2015. Innovative diagnostic tools for early detection of Alzheimer's disease. *Alzheimer's & Dementia* 11(5), 561-578.
- Laufer B., Nation, P. 1995. Vocabulary Size and Use: Lexical Richness in L2 Written Production. *Applied Linguistics*. 16:3. pp. 307-322.
- Masrani, V., Murray, G., Shoshana Field, T., Carenini, G. 2017. Domain Adaptation for Detecting MCI. *Advances in AI*. Pp 248-259. Springer.
- McCarthy, P., Jarvis, S. 2010. MTLT, voc-d, and HD-d: *Beh Res Meth.* 42, 381-392.
- Milton, J. 2009. *Measuring Second Language Vocab Acquisition*. *Multiling Matters*.
- Nation, I.S.P. 2013. *Learning Vocabulary in Another Language*. *Cambr. Appl. Ling.*
- Prince, M. et al. 2013. The global prevalence of dementia. *Alz & Dem* 9(1), 63-75.
- Richards B., Malvern, D.D., Meara, P., Milton, J., Treffers-Daller, J. 2009. *Vocabulary Studies in 1st and 2nd Language Acquisition. The Interface Between Theory and Application*. Palgrave Macmillan.
- Wallin A., Nordlund A., Jonsson M., Lind K., Edman Å., Göthlin M., Stålhammar J., Eckerström M., Kern S., Börjesson-Hanson A., Carlsson M., Olsson E., Zetterberg H., Blennow K., Svensson J., Öhrfelt A., Bjerke M., Rolstad S., Eckerström C. The Gothenburg mci study. *J of Cerebral Blood Flow and Metabolism*, 36(1):114–31.