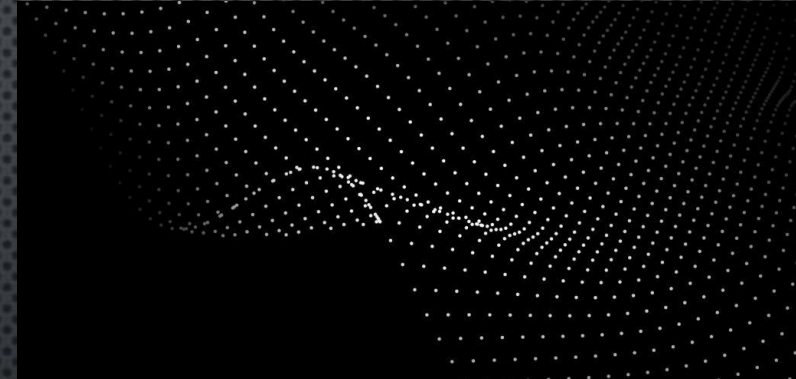


# PREDICTIVE AND DISTINCTIVE LINGUISTIC FEATURES IN SCHIZOPHRENIA-BIPOLAR SPECTRUM DISORDERS



MARTINA KATALIN SZABÓ, VERONIKA VINCZE, BERNADETT DAM, CSENGE GUBA, ANITA BAGI, ISTVÁN SZENDI

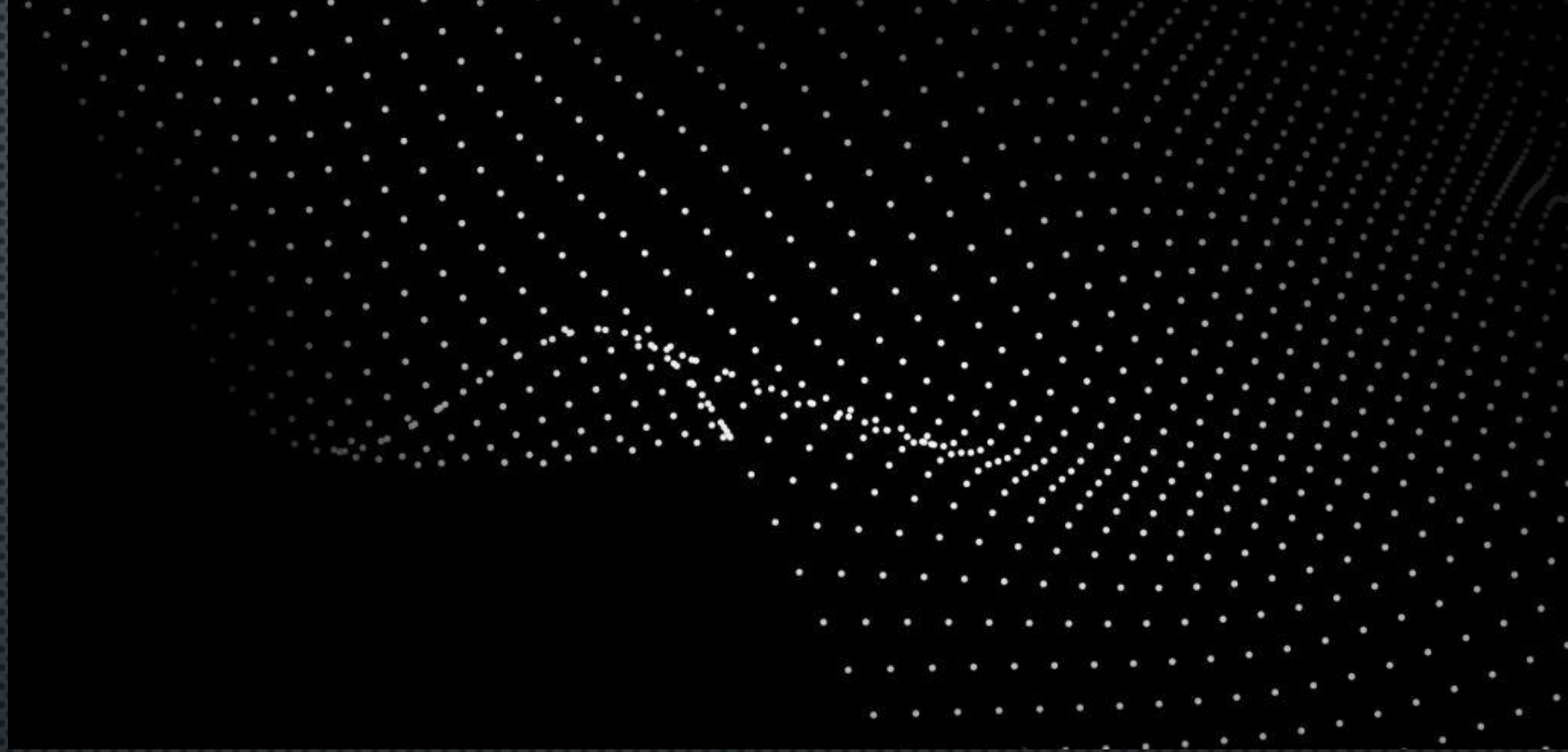
University of Szeged, Hungary; HUN-REN Centre for Social Sciences, Hungary; Tokyo University of Foreign Studies, Japan; HUN-REN-SZTE Research Group on Artificial Intelligence, Hungary; Kiskunhalas Semmelweis Hospital, Teaching Hospital of the University of Szeged Psychiatry Department, Hungary

Szabo.Martina@tk.hu, vinczev@inf.u-szeged.hu, dam.bernadett@stud.u-szeged.hu, csenge.guba@gmail.com, bagianita88@gmail.com, szendi@inf.u-szeged.hu



# Introduction

- Spontaneous speech analysis of Hungarian patients with schizophrenia, schizoaffective, and bipolar disorders
- Comprehensive statistical analysis of linguistic parameters
- Goal: Identify distinctive linguistic features
- Method: Automatic differentiation among patient groups and controls using random forest algorithm
- Results: Effective distinction among SZ, SAD, BD, and controls, surpassing baseline results.



Understanding the relationship between  
psychosis spectrum disorders and  
language behavior



# Understanding the relationship between psychosis spectrum disorders and language behavior

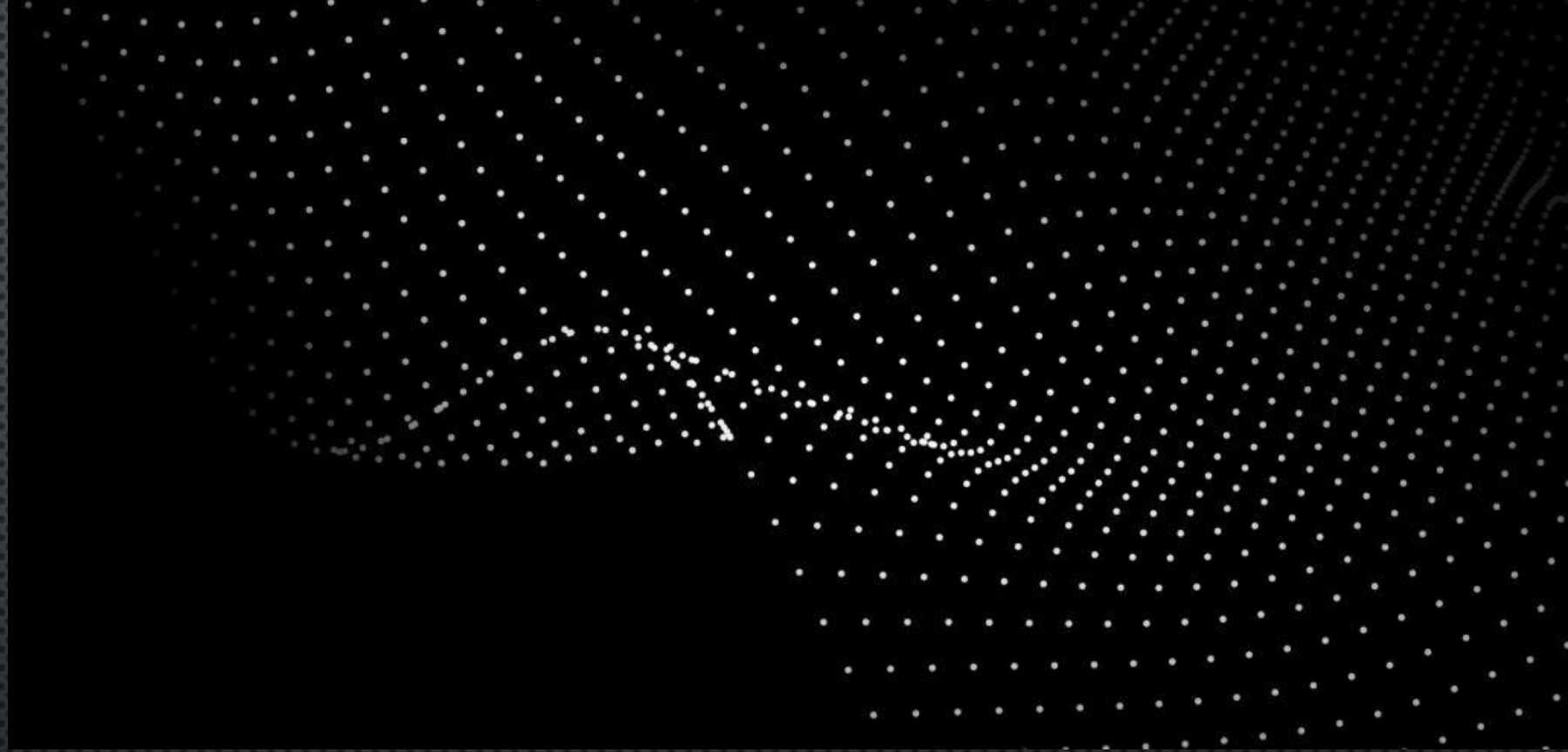
- Bipolar disorder (BD)
  - Recurrent episodes of mania, hypomania, and depression.
- Schizophrenia (SZ)
  - Symptoms: delusions, hallucinations, disorganized behavior.
- Schizoaffective disorder (SAD)
  - Mixed psychotic and affective symptoms.
- Cognitive impairment
  - Common in psychotic disorders.
  - May influence language use.
- Mental health and communication
  - Analysis of linguistic data provides insights into the relationship between linguistic factors and psychological aspects.

# Research Gap



- **Limited Analysis:**
  - Only one research work compares linguistic features of text produced by SZ, SAD, and BD (Voleti et al., 2019).
  - However, this analysis does not distinguish between patients with schizophrenia or schizoaffective disorder.
- **Automatic Discrimination:**
  - Automatic discrimination among these patient groups based on linguistic features has not yet been addressed in the literature.
- **Hungarian Population Studies:**
  - Several papers have examined Hungarian patients with SZ, SAD, and BD (e.g., Kéri et al., 2001; Réthelyi et al., 2010; Inczédy et al., 2010; Kocsis et al., 2016; Döme et al., 2005; Kárpáti et al., 2018).
  - However, no study has systematically analyzed the linguistic features of these disorders in the Hungarian population.





Corpus compilation

# Data collection

- Recorded by the Prevention of Mental Illnesses Interdisciplinary Research Group, University of Szeged, led by István Szendi
- Data collection approved by the Ethics Committee of the University of Szeged, following the Declaration of Helsinki
- Written informed consent and official permission obtained
- Spontaneous speech reflects language specificities more accurately than planned speech
- The database includes 458 monologues from 77 subjects

	Groups				
	Control	SZ	SAD	BD	all
Number of participants	21	27	14	15	77
Number of texts	126	162	84	86	458
Age; M(SD)	36.42(10.49)	38.80(10.17)	41.43(9.73)	49.08(8.67)	40.63(10.71)
Education; M(SD)	14.76(3.05)	14.19(2.87)	14.54(2.98)	15.83(3.62)	14.73(3.09)
Sex ratio; f:m	8:13	9:18	10:4	8:7	35:42

Table 1: Basic Data of the Four Subject Groups. Education is given in years.  
(M = mean, SD = standard deviation)

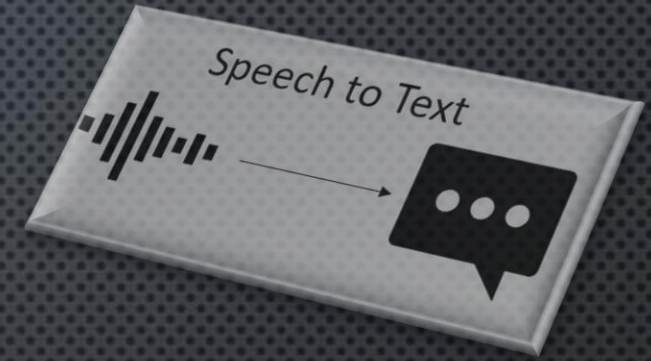
## Speech tasks:

Description of Self (DescSelf), Mother (DescMother), and Father (DescFather), Younger Years (YoungSelf and YoungOther), Description of Previous Day (PrevDay)

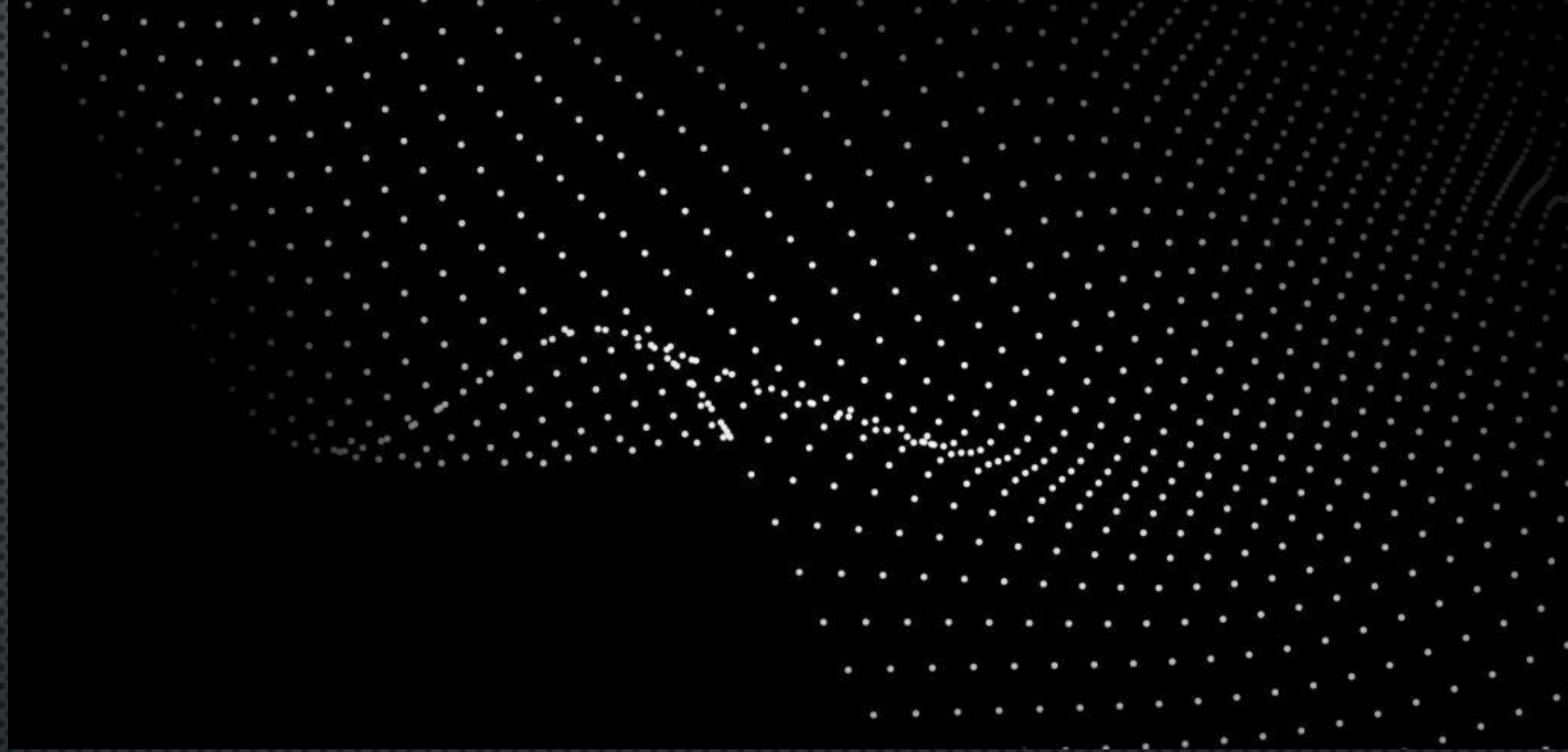


# Transcription process

- Manual transcription of recordings
- No specific software used
- Marked pauses, hesitations etc.
- Transcriptions stored in plain text format (UTF-8)
- Transcribers worked interchangeably on files
- Detailed guideline used for consistency
- Regular quality checks conducted







Corpus processing and analysis

# Corpus processing

- We used the magyarlanc toolkit for automatic linguistic analysis, including sentence splitting, tokenization, lemmatization, part-of-speech tagging, and morphological tagging.
- Extracted: 17 basic statistical features and 10 speech-based, 87 morphosyntactic features
- Semantic and pragmatic linguistic features were analyzed, including sentiment and emotion words, discourse markers, and intensifiers.
- A total of 194 linguistic features

## Example Features:

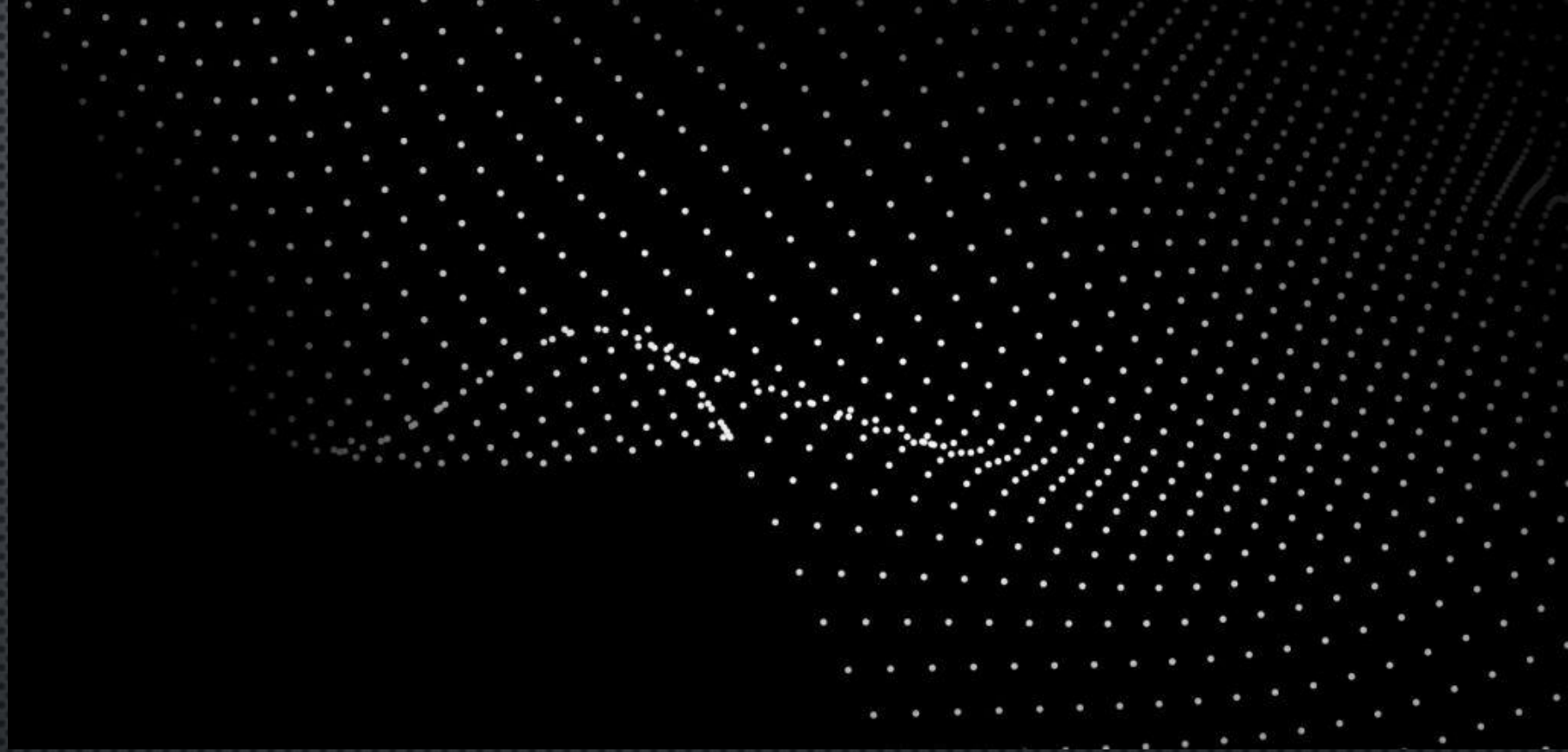
- Basic Statistical: Number of sentences, number and frequency of distinct lemmas compared to the number of words.
- Speech-Based: Number of pauses, hesitations, and filled pauses.
- Morphosyntactic: Count and occurrence rate of various parts-of-speech, and occurrences of third person singular verb forms.
- Semantic and Pragmatic: Number and frequency of sentiment words, emotion words, and discourse markers.

(For a detailed list of features, please refer to the appendix of our paper.)



# Statistical analysis and machine learning experiments

- Pairwise t-tests -> evaluate feature usefulness
- Random forest classifier from the WEKA package -> automatic group discrimination
- Ten-fold cross-validation
- Results compared to majority classification baseline method



Corpus analysis results

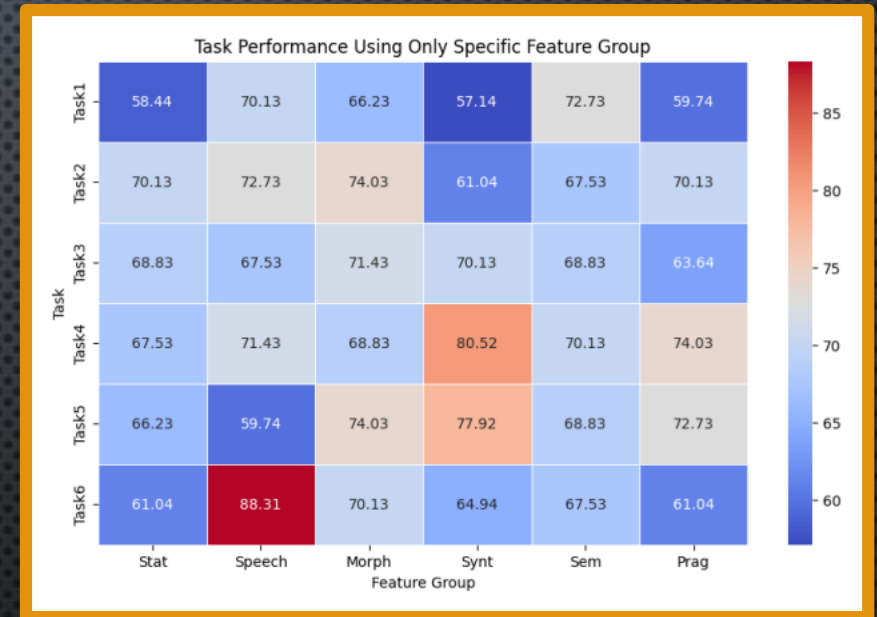


# Results of significance tests

- Several significant differences among the four groups:
  - The rate of **nouns and verbs** was significantly different in most sub-corpora.
  - Patients tended to use fewer **pronouns** compared to the control group.
  - Patients tended to use fewer **function** and more **content words** generally.
- When comparing each patient group with the control group, clear patterns of significant differences emerged:
  - The most apparent contrast was seen in BD patients, who had a lower rate of **filled and unfilled pauses** in several tasks.
  - SZ patients had a notably different **POS** distribution compared to the control group.
  - Patients used significantly more **positive words** when discussing someone close to them, but more **negative words** when discussing themselves.
  - BD speakers also used significantly more words related to **sorrow**
  - Healthy controls typically used the most **discourse markers** across all tasks.

# Results of machine learning experiments I.

- Overall accuracy: 58.44%, outperforming the baseline.
- Dividing subjects into two groups yielded an accuracy of 72.73%
- Best performance: SZ subjects (F-score: 0.75).
- Speech-based features were the most effective (63.64%).

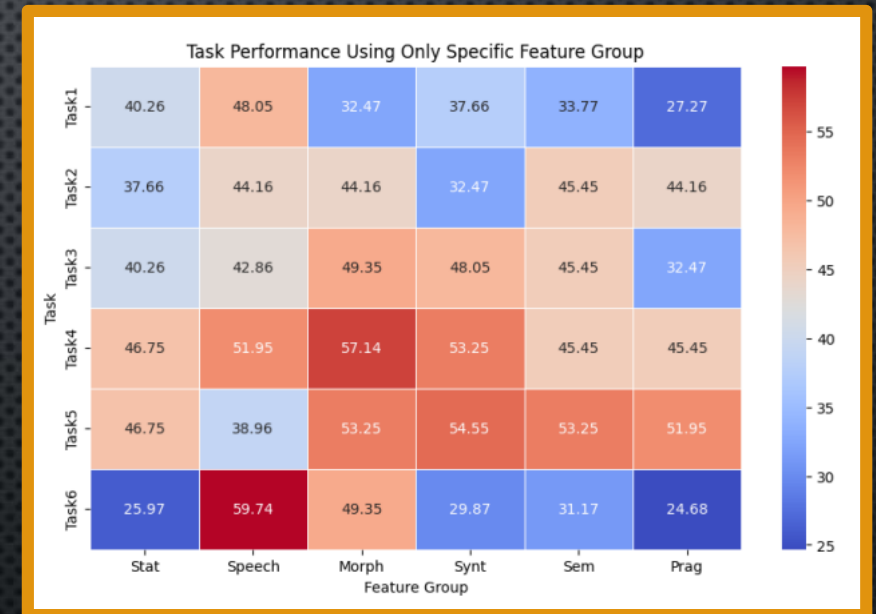


- Highest efficiencies by task and feature:
  - "PrevDay" with speech-based features (88.31%).
  - "YoungSelf" with syntactic features (80.52%).
  - "YoungOther" with syntactic features (77.92%).



# Results of machine learning experiments II.

- When comparing the four groups:  
Excluding statistical or syntactic features improved accuracy
- Omitting some specific morphological and semantic features increased accuracy



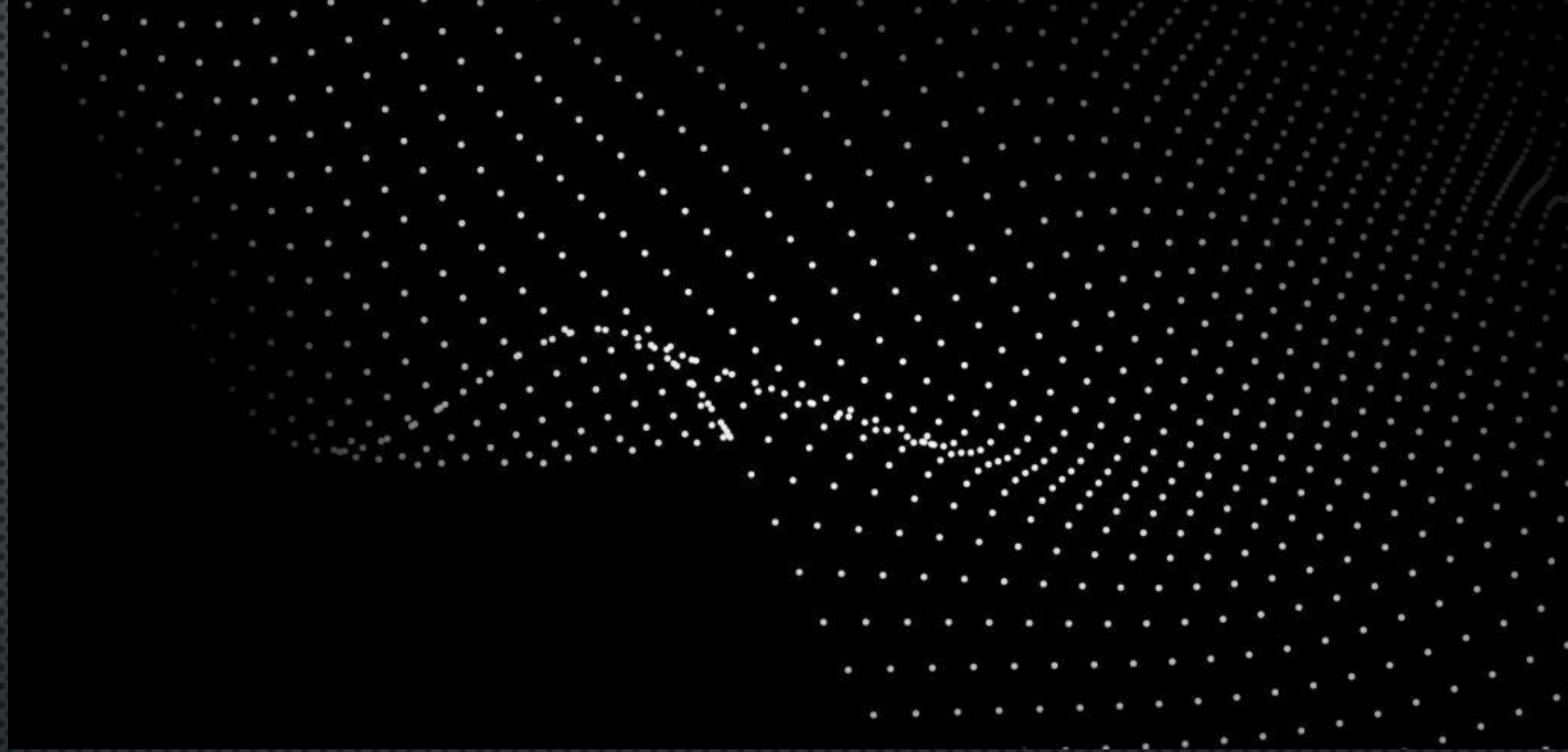
# Results of machine learning experiments III.

- In our analysis and machine learning experiments, we also conducted experiments with intensifiers.
- Intensifiers are closely linked to emotion regulation (Athanasiadou, 2007; Strous et al., 2009); There are differences in the use of intensifiers among various mental disorders, such as schizophrenia (Strous et al., 2009, Szabó et al. 2023).
- Two types: Standard-register (non-emotive) intensifiers and Negative emotive intensifiers (NEI) (e.g. *awfully*, *crazy*, *damn*, *brutally* etc.)
- **Machine Learning Experiments:**
  - Excluding NEIs reduced overall accuracy but removing standard-register intensifiers increased accuracy
  - Including NEIs improved machine learning results

## Importance of Feature Selection:

- These findings underscore the significance of feature selection.
- They highlight the complex interaction among different feature groups, significantly enhancing automatic classification performance.





Conclusions and future work

# Conclusion and future work

- In this study, we analyzed spontaneous speech from Hungarian patients with SZ, SAD, and BD, aiming to identify distinctive linguistic features and improve automatic classification.
- Overall accuracy: 58.44%; Controls from patients: accuracy of 72.73%

## **Future work:**

- Expand linguistic feature set; conduct deeper analysis
- Fine-tune machine learning methods
- Focus on most effective tasks and features
- Experiment with various machine learning algorithms
- Investigate data from other languages
- Release the corpus for research purposes after masking sensitive data.



## REFERENCES

- ANGELIKI ATHANASIADOU. 2007. ON THE SUBJECTIVITY OF INTENSIFIERS. *LANGUAGE SCIENCES*, 29(4):554–565.
- P DÖME, Z RIHMER, X GONDA, P PESTALITY, G KOVÁCS, Z TELEKI, AND P MANDL. 2005. CIGARETTE SMOKING AND PSYCHIATRIC DISORDERS IN HUNGARY. *INTERNATIONAL JOURNAL OF PSYCHIATRY IN CLINICAL PRACTICE*, 9(2):145–148.
- GABRIELLA INCZÉDY-FARKAS, JUDIT BENKOVITS, NÓRA BALOGH, PÉTER ÁLMOS, BEÁTA SCHOLTZ, GÁBOR ZAHUCZKY, ZSOLT TÖRÖK, KRISZTIÁN NAGY, JÁNOS RÉTHELYI, ZOLTÁN MAKKOS, ET AL. 2010. SCHIZOBANK – THE HUNGARIAN NATIONAL SCHIZOPHRENIA BIOBANK AND ITS ROLE IN SCHIZOPHRENIA RESEARCH. *ORVOSI HETILAP*, 151(35):1403–1408.
- ESZTER KÁRPÁTI, ANITA BAGI, ISTVÁN SZENDI, LUJZA BEATRIX TÓTH, KAROLINA JANACSEK, AND ILDIKÓ HOFFMANN. 2018. REKURZÍÓ EGY SZKIZOAFFEKTÍV ZAVARRAL ÉLŐ SZEMÉLY DISKURZUSAIBAN– ESETTANULMÁNY. *ISKOLAKULTÚRA: PEDAGÓGUSOK SZAKMAI-TUDOMÁNYOS FOLYÓIRATA*, 28(5- 6):40–54.
- SZABOLCS KÉRI, O KELEMEN, G BENEDEK, AND Z JANKA. 2001. DIFFERENT TRAIT MARKERS FOR SCHIZOPHRENIA AND BIPOLAR DISORDER: A NEUROCOGNITIVE APPROACH. *PSYCHOLOGICAL MEDICINE*, 31(5):915–922.
- KRISZTINA KOCSIS-BOGÁR, ZSÓFIA NEMES, AND DÓRA PERCZEL-FORINTOS. 2016. FACTORIAL STRUCTURE OF THE HUNGARIAN VERSION OF OXFORD-LIVERPOOL INVENTORY OF FEELINGS AND EXPERIENCES AND ITS APPLICABILITY ON THE SCHIZOPHRENIA-SCHIZOTYPY CONTINUUM. *PERSONALITY AND INDIVIDUAL DIFFERENCES*, 90:130–136.
- JÁNOS M RÉTHELYI, STEVEN C BAKKER, PATRÍCIA POLGÁR, PÁL CZOBOR, ERIC STRENGMAN, PÉTER I PÁSZTOR, RENÉ S KAHN, AND ISTVÁN BITTER. 2010. ASSOCIATION STUDY OF NRG1, DTNBP1, RGS4, G72/G30, AND PIP5K2A WITH SCHIZOPHRENIA AND SYMPTOM SEVERITY IN A HUNGARIAN SAMPLE. *AMERICAN JOURNAL OF MEDICAL GENETICS PART B: NEUROPSYCHIATRIC GENETICS*, 153(3):792–801.
- MARTINA KATALIN SZABÓ, VERONIKA VINCZE, CSERGE GUBA, BERNADETT DAM, ADRIENN SOLYMOS, ANITA BAGI, AND ISTVÁN SZENDI. 2023. FOKOZÁS SZKIZOFRÉNIÁBAN. IN XIX. MAGYAR SZÁMÍTÓGÉPES NYELVÉSZETI KONFERENCIA, PAGES 17–32.
- ROHIT VOLETI, STEPHANIE WOOLRIDGE, JULIE M LISS, MELISSA MILANOVIC, CHRISTOPHER R BOWIE, AND VISAR BERISHA. 2019. OBJECTIVE ASSESSMENT OF SOCIAL SKILLS USING AUTOMATED LANGUAGE ANALYSIS FOR IDENTIFICATION OF SCHIZOPHRENIA AND BIPOLAR DISORDER. *ARXIV PREPRINT ARXIV:1904.10622*.
- RAE D STROUS, MOSHE KOPPEL, JONATHAN FINE, SMADAR NACHLIEL, GINETTE SHAKED, AND ARI Z ZIVOTOFSKY. 2009. AUTOMATED CHARACTERIZATION AND IDENTIFICATION OF SCHIZOPHRENIA IN WRITING. *THE JOURNAL OF NERVOUS AND MENTAL DISEASE*, 197(8):585–588.

## ACKNOWLEDGE MENTS

THIS WORK WAS SUPPORTED BY THE EFOP-3.6.1-16-2016-00008 (SUB-PROJECT: „SZKIZOFRÉRIA – NYELV – INNOVATÍV GYÓGYÍTÁS”), HUN-REN CENTRE FOR SOCIAL SCIENCES, THE UNIVERSITY OF SZEGED, DEPARTMENT OF INFORMATICS AND THE INSTITUTE OF GLOBAL STUDIES, TOKYO UNIVERSITY OF FOREIGN STUDIES, JAPAN. \\ MARTINA KATALIN SZABÓ WAS FUNDED BY THE OTKA POSTDOCTORAL RESEARCH GRANT, THE HUNGARIAN RESEARCH FUND OF THE NATIONAL RESEARCH, DEVELOPMENT AND INNOVATION OFFICE OF HUNGARY (NKFIH) (GRANT NUMBER: PD 132312) AND BY THE INTERNATIONAL RESEARCH FELLOWSHIP PROGRAM OF JAPAN SOCIETY FOR THE PROMOTION OF SCIENCE (JSPS, POSTDOCTORAL FELLOWSHIPS FOR RESEARCH IN JAPAN (STANDARD)).





THANK YOU FOR YOUR  
ATTENTION.