

**Mental Health Evaluation through Text Analysis (META)**

Gabriel Bonnin

Ruhr University Bochum, Germany

**Author Note**

Correspondence concerning this article should be addressed to Gabriel Bonnin, Email:  
[gabriel.bonnin@ruhr-uni-bochum.de](mailto:gabriel.bonnin@ruhr-uni-bochum.de)

## Mental Health Evaluation through Text Analysis (META)

### Abstract

Psychotherapy is one of the most effective treatments for mental health problems, but its success depends on accurate diagnostic assessments. Current diagnostic tools often use standardized closed-ended scales that, while reliable, may fail to capture the complexity and individuality of mental states. In collaboration with Dr. Oscar Kjell at Lund University, this project leverages advances in artificial intelligence (AI) and natural language processing (NLP) to transform the way mental health is assessed.

This project is based on a unique longitudinal dataset collected over 10 years at the Mental Health Research and Treatment Center (FBZ) at Ruhr University Bochum. It consists of written self-reports in which patients describe their mental health problems, functional impairment, and therapy goals in their own words. These texts are linked to key clinical outcomes such as diagnoses, symptom severity, functional impairment, and treatment success, providing a rich, ecologically valid resource for understanding patient progress in psychotherapy.

While previous NLP-based mental health studies have focused on social media posts, limiting their clinical applicability, this project applies state-of-the-art large language models (LLMs) to real-world clinical data. By analyzing patients' open-ended responses, we aim to uncover patterns in how they articulate mental health, emotions, and treatment trajectories. These insights will inform the development of AI-powered tools that offer more personalized and clinically relevant assessments, surpassing traditional methods in accuracy and depth. Ultimately, this research aims to support clinicians in making more informed treatment decisions, enhance personalized care, and contribute to the modernization of mental health assessment.

### Introduction

Mental health problems pose a significant global challenge, accounting for a considerable proportion of deaths and disability-adjusted life years ([World Health Organization, 2017](#)). Psychotherapy is an effective and sustainable intervention for reducing symptoms and improving quality of life ([Chorpita et al., 2011](#); [Wampold & Imel, 2015](#)), but

its success critically depends on accurate assessments (Jensen-Doss & Weisz, 2008; Lutz et al., 2022).

Current assessment practices typically combine subjective self-reports with clinical observations (Bonnin et al., 2024). Standardized closed-ended tools such as the Beck Depression Inventory-II (Beck et al., 1996) rely on numerical scales (Likert, 1932) to structure and standardize assessments and are widely used in clinical research and practice. While these methods have advanced replicability and reliability in psychological assessment, they can miss important individual differences by restricting responses to pre-defined categories, limiting the ability to capture the complexity of mental states (Kjell, Kjell, et al., 2024).

Recent advances in AI, particularly transformer-based LLMs (Vaswani et al., 2017), present promising solutions to these limitations (Kjell, Kjell, et al., 2024). LLMs excel in analyzing context-rich natural language with remarkable accuracy across diverse tasks (Devlin et al., 2019). Open-ended response formats, where patients describe their experiences in their own words, provide high-dimensional, context-rich information that remains underutilized in current assessment practices. Empirical studies highlight the potential of NLP-based analysis of open-ended responses, achieving moderate convergence with closed-ended rating scales using traditional NLP methods (Kjell et al., 2019) and nearing theoretical upper limits of accuracy with LLMs (Kjell et al., 2022). Preliminary research also highlights their potential for predicting clinically significant outcomes, including suicide risk (Matero et al., 2019; Mohammadi et al., 2019; Zirikly et al., 2019).

At the FBZ at Ruhr University Bochum, open-ended patient responses have been routinely collected from approximately 3,000 patients pre-therapy. While previous studies have explored the use of NLP for mental health assessment, they have largely relied on social media language, limiting their clinical applicability. My project moves beyond the current state-of-the-art by applying LLMs to analyze this unique, large-scale, and longitudinal clinical dataset, assessing the relationship between patients' probed mental health responses and clinically relevant constructs such as diagnosis, symptom severity, and functional impairment. Additionally, it will predict key clinical outcomes such as treatment response and therapy goal attainment. Furthermore, the project seeks to generate clinically meaningful, data-driven

insights that go beyond traditional diagnostic categories, offering a more nuanced and patient-centered understanding of mental health trajectories.

## Methods

### Preprocessing

To streamline data collection, an automated transcription pipeline was implemented: The handwritten text data is first recorded by trained employees of the FBZ adult outpatient clinic using a mobile audio recording device. Identifying features (e.g. names, dates of birth, location details) were replaced by placeholders during recording (anonymisation). The transcription was carried out on local hardware using the open source tool Whisper Large v2 (<https://github.com/openai/whisper>), a state-of-the-art speech-to-text model (Radford et al., 2022). Each recording begins with a structured introduction, including a patient identification code, followed by responses to predefined questions. The transcription pipeline automatically processes all audio recordings, extracts the patient codes, and identifies responses to key questions.

As an additional data correction step, the exported transcription table was screened for incomplete entries. Records with missing patient identification codes or without any extracted text were automatically flagged, exported for manual correction, and subsequently re-imported and merged back into the original dataset. The corrected dataset was then used for downstream analyses.

### Measures

#### *Sociodemographic and context measures*

Sociodemographic information included age, sex, marital and relationship status, general educational attainment, vocational qualification, and current work ability. Contextual variables captured prior psychological or psychiatric treatment and the manner in which therapy ended (e.g., regular completion, dropout).

#### *Responses from open-ended questions before therapy*

At the start of therapy, patients complete two separate questionnaires designed to assess key aspects of their mental health concerns, functional impairments, and expectations

for treatment. Questions 1–9 come from the first questionnaire (*Fragebogen zur Lebensgeschichte*), and questions 10–13 come from the second (*Eingangsfragebogen*). The questions include:

1. **Problem development:** ‘Briefly describe how the problems for which you are seeking treatment have developed over time.’ (geman original question: „Beschreiben Sie kurz, wie sich Ihre Probleme, wegen derer Sie eine Behandlung aufsuchen, im Laufe der Zeit entwickelt haben.“)
2. **Extra stressors:** ‘What causes you stress in addition to your everyday problems (e.g. finances, housing situation)?’ (geman original question: „Was macht Ihnen zusätzlich zu Ihren Problemen im Alltag Stress (z. B. Finanzen, Wohnsituation);“)
3. **Pre-onset changes:** ‘Did something special change in your life before the onset of your symptoms? (e.g. death of an important person, divorce or separation, change in work situation or income, addition to the family)’ (geman original question: „Hat sich vor dem Beginn Ihrer Beschwerden etwas Besonderes in Ihrem Leben verändert? (z. B. Tod einer wichtigen Bezugsperson, Scheidung oder Trennung, Veränderung der Arbeitssituation oder des Einkommens, Familienzuwachs)“)
4. **Event connection:** ‘Do you see a connection between the event(s) and the development of your problems?’ (geman original question: „Sehen Sie einen Zusammenhang zwischen dem Ereignis/den Ereignissen und der Entwicklung Ihrer Probleme;“)
5. **Physical symptoms:** ‘Are there any physical side effects when your problems occur?’ (geman original question: „Gibt es körperliche Begleiterscheinungen, wenn Ihre Probleme auftreten;“)
6. **Problem causes:** ‘What do you think are the causes of your problems?’ (geman original question: „Welche Ursachen sehen Sie für Ihre Probleme;“)
7. **Expected improvements:** ‘What would improve in your life if you no longer had your problems?’ (geman original question: „Was würde sich in Ihrem Leben verbessern, wenn Sie ihre Probleme nicht mehr hätten;“)
8. **Environment response:** ‘How does your environment (partner, family, friends, work colleagues) react to your problems?’ (geman original question: „Wie reagiert Ihre

- Umwelt (Partner:in, Familie, Freund:innen, Arbeitskolleg:innen) auf die Probleme;‘)
9. **No change required:** ‘What should not change under any circumstances as a result of the therapy?’ (geman original question: „Was sollte sich durch die Therapie auf keinen Fall verändern;‘)
  10. **Problem description:** ‘Finally, please describe in your own words the problems for which you would like treatment.’ (geman original question: „Beschreiben Sie zum Abschluss bitte noch einmal in eigenen Worten Ihre Probleme, deretwegen Sie eine Behandlung wünschen.“)
  11. **Impacted life areas:** ‘In which areas of your life do these problems limit you (e.g. job, relationship)?’ (geman original question: „In welchen Lebensbereichen schränken Sie diese Probleme ein (z. B. Beruf, Partnerschaft);‘)
  12. **Therapy goals:** ‘What would you like to achieve for yourself in therapy?’ (geman original question: „Was möchten Sie in der Therapie für sich erreichen;‘)

### ***Psychometric measures***

Clinical and psychometric variables were retrieved from the FBZ database and included diagnostic information, self-report symptom measures, therapist- and patient-rated outcome measures, positive mental health indicators, and therapeutic process variables. Diagnoses were coded according to DSM-5 and ICD-10 criteria. Symptom severity and treatment outcomes were assessed using a combination of standardized self-report questionnaires and clinician-rated instruments administered at different points during treatment.

**Diagnosis.** Diagnosis at the outpatient clinic is conducted using structured clinical interviews. These typically take place before therapy begins, usually at the fourth therapist–patient contact. The interview used is the Diagnostic Interview for Mental Disorders ([Margraf et al., 2021](#)), which covers the most frequent DSM-5 disorders encountered in outpatient therapy settings.

**Beck-Depression-Inventory II..** Depressive symptoms were assessed using the *Beck Depression Inventory-II* (BDI-II; ([Beck et al., 1996](#))), a widely used self-report questionnaire measuring the severity of depressive symptomatology over the past two weeks.

**Depression Anxiety Stress Scale 42.** General psychological distress was assessed using the *Depression Anxiety Stress Scales–42* (DASS-42([Lovibond & Lovibond, 1995](#))), which consists of 42 items measuring symptoms of depression, anxiety, and stress on a XXX scale.

**Brief Symptom Inventory.** Overall psychopathological symptom burden was measured using the *Brief Symptom Inventory* (BSI; ([Franke, 2002](#))), the short form of the Symptom Checklist-90-Revised (SCL-90-R; Derogatis). The BSI consists of 53 items rated on a 5-point Likert scale ranging from 0 (“not at all”) to 4 (“extremely”). Responses to 49 items are assigned to nine primary symptom dimensions, while four items are evaluated separately. These symptom dimensions are summarized into three global indices: the *Global Severity Index* (GSI), reflecting overall psychological distress; the *Positive Symptom Distress Index* (PSDI), indicating symptom intensity; and the *Positive Symptom Total* (PST), representing the number of reported symptoms.

**Positive Mental Health Scale.** Positive mental health (PMH) was assessed with the nine-item PMH scale ([Lukat et al., 2016](#)). Responses are given on a 4-point Likert scale from 0 (disagree) to 3 (agree). Item scores are summed to yield a total score ranging from 0 to 27, with higher scores reflecting greater PMH. The scale has been validated as a unidimensional measure with excellent internal consistency (Cronbach’s  $\alpha = .93$ ), good test-retest reliability (Pearson  $r = .74\text{--}.81$ ), and evidence of scalar invariance across samples and over time ([Lukat et al., 2016](#)). Furthermore, it shows strong convergent and discriminant validity and is sensitive to therapeutic change across diverse populations ([Lukat et al., 2016](#)).

**Childhood Trauma Questionnaire.** Early adverse experiences were assessed using the *Childhood Trauma Questionnaire* (CTQ), a widely used self-report instrument for the retrospective assessment of childhood maltreatment. The CTQ measures five domains of adverse experiences: emotional abuse, physical abuse, sexual abuse, emotional neglect, and physical neglect. Items are rated on a 5-point Likert scale ranging from 1 (“not at all”) to 5 (“very often”), with higher scores indicating greater exposure to maltreatment. The German version of the CTQ has demonstrated good psychometric properties, including satisfactory reliability and validity in clinical samples

Wingenfeld, Katja & Spitzer, Carsten & Mensebach, Christoph & Grabe, Hans & Hill, Andreas & Gast, Ursula & Schlosser, Nicole & Höpp, Hella & Beblo, Thomas & Driessen, Martin. (2010). The German Version of the Childhood Trauma Questionnaire (CTQ): Preliminary Psychometric Properties. Psychotherapie, Psychosomatik, medizinische Psychologie. 60. 10.1055/s-0030-1253494.

### **Clinical Global Impression.**

**Severity Scale.** Clinician-rated symptom severity and improvement were assessed using the *Clinical Global Impression* (CGI) scales. The *CGI-Severity* scale evaluates the clinician's global impression of the patient's current level of mental illness, based on their total clinical experience with this population. The item asks: "Considering your total clinical experience with this particular population, how mentally ill is the patient at this time?"

**Improvement Scale.** Treatment-related change was assessed using the *CGI-Improvement* scale. Both patients and therapists rated overall improvement relative to the beginning of therapy, regardless of whether the change was attributed entirely to treatment. Patient and therapist versions differ only in perspective but use equivalent response formats.

### **Global Improvement**

Global therapy outcome was assessed using a six-point global success rating based on two items measuring perceived benefit and satisfaction with therapy ([Michalak et al., 2003](#)). These items were completed by both patients and therapists and capture a retrospective evaluation of treatment success. The items assess (1) the extent to which expectations toward therapy have been fulfilled and (2) the overall perceived benefit of therapy. Responses are given on a 6-point Likert scale ranging from 1 ("on the contrary / rather harmful") to 6 ("completely / very helpful").

### **Goal Attainment Scale**

Individualized treatment outcomes were assessed using a goal attainment measure inspired by the Goal Attainment Scaling approach ([Kiresuk & Sherman, 1968](#)). At the beginning of therapy, patients and therapists collaboratively define individualized treatment goals. During the course and at the end of therapy, patients and therapists retrospectively evaluated the extent to which each of the predefined goals had been achieved.

Goal attainment was rated on a standardized six-point numerical scale ranging from deterioration relative to the initial goal state (-1 = moved away from the goal) to full goal attainment (4 = goal achieved). The scale reflects patients' subjective assessment of goal progress, with intermediate categories indicating partial progress toward the respective goal.

For each patient, an overall goal attainment score was computed as the mean rating across all individually defined goals, representing the average subjective level of goal progress at the end of therapy.

### ***Measurement time points***

### **Analysis**

The Sequential Evaluation with Model Pre-registration (Kjell, Ganesan, et al., 2024) framework will be implemented to ensure robust model development, mitigating overfitting and enabling unbiased performance evaluation. Additionally, evaluating models on prospective data will simulate real-world clinical deployment by assessing performance on new, unseen patient data.

During the model development phase, preprocessing pipelines will be finalized, and exploratory models developed using advanced cross-validation techniques. Contextual embeddings derived from pretrained LLMs will be linked to clinical outcomes using state-of-the-art prediction models, including ridge regression (Hoerl & Kennard, 1970), lasso regression (Tibshirani, 1996), and random forests (Ho, 1995). The final pipelines will be pre-registered for evaluation (e.g., <https://aspredicted.org/>).

In the evaluation phase, pre-registered models will be tested on held-out datasets, enabling unbiased performance assessments. This phase will include detailed error analysis to identify potential biases and limitations, as well as comparative analyses to benchmark model performance against the HRG's models.

Furthermore, topic modeling techniques (e.g., Latent Dirichlet Allocation; (Blei et al., 2003), or BERTopic; (Grootendorst, 2022)) will be employed to explore themes in patient-generated text. These analyses will provide valuable clinical insights into patients' subjective experiences, including their perceived problems, impairments, and goals. By analyzing these insights, I aim to highlight commonalities and differences in patient narratives

across diverse populations or conditions. This process may also reveal nuanced linguistic cues that correlate with clinical outcomes, offering a richer understanding of patient perspectives and informing personalized care strategies.

## Results

## References

- Beck, A. T., Steer, R. A., & Brown, G. K. (1996). *BDI-II: Beck Depression Inventory manual* (2nd ed.). Psychological Corporation.
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3(null), 993–1022.
- Bonnin, G., Kröber, S., & Brachel, R. von. (2024). Die Klassifikation psychischer Störungen mit diagnostischen Interviews. In T. Teismann, P. Thoma, S. Taubner, A. Wannemüller, & K. von Sydow (Eds.), *Klinische Psychologie und Psychotherapie: Ein verfahrensübergreifendes Lehr- und Lernbuch*. Hogrefe.
- Chorpita, B. F., Daleiden, E. L., Ebetsutani, C., Young, J., Becker, K. D., Nakamura, B. J., Phillips, L., Ward, A., Lynch, R., Trent, L., et al. (2011). Evidence-based treatments for children and adolescents: An updated review of indicators of efficacy and effectiveness. *Clinical Psychology: Science and Practice*, 18(2), 154–172.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), *Proceedings of the 2019 conference of the north American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers)* (pp. 4171–4186). Association for Computational Linguistics.  
<https://doi.org/10.18653/v1/N19-1423>
- Franke, G. (2002). *Franke, g.h. (2000). BSI. Brief symptom inventory - deutsche version. Manual*. Göttingen: beltz.
- Grootendorst, M. (2022). *BERTopic: Neural topic modeling with a class-based TF-IDF procedure*. <https://arxiv.org/abs/2203.05794>
- Ho, T. K. (1995). Random decision forests. *Proceedings of 3rd International Conference on Document Analysis and Recognition*, 1, 278–282 vol.1.

<https://doi.org/10.1109/ICDAR.1995.598994>

Hoerl, A. E., & Kennard, R. W. (1970). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics*, 12(1), 55–67.

<https://doi.org/10.1080/00401706.1970.10488634>

Jensen-Doss, A., & Weisz, J. R. (2008). Diagnostic agreement predicts treatment process and outcomes in youth mental health clinics. *Journal of Consulting and Clinical Psychology*, 76(5), 711–722. <https://doi.org/10.1037/0022-006X.76.5.711>

Kiresuk, T. J., & Sherman, R. E. (1968). Goal attainment scaling: A general method for evaluating comprehensive community mental health programs. *Community Mental Health Journal*, 4(6), 443–453. <https://doi.org/10.1007/BF01530764>

Kjell, O. N. E., Ganesan, A. V., Boyd, R., Oltmanns, J. R., Rivero, A., Feltman, S., Carr, M. A., Luft, B. J., Kotov, R., & Schwartz, H. A. (2024). *Demonstrating high validity of a new AI-language assessment of PTSD: A sequential evaluation with model pre-registration*. PsyArXiv. <https://doi.org/10.31234/osf.io/xw24e>

Kjell, O. N. E., Kjell, K., Garcia, D., & Sikström, S. (2019). Semantic measures: Using natural language processing to measure, differentiate, and describe psychological constructs. *Psychological Methods*, 24(1), 92–115. <https://doi.org/10.1037/met0000191>

Kjell, O. N. E., Kjell, K., & Schwartz, H. A. (2024). Beyond rating scales: With targeted evaluation, large language models are poised for psychological assessment. *Psychiatry Research*, 333, 115667. <https://doi.org/10.1016/j.psychres.2023.115667>

Kjell, O. N. E., Sikström, S., Kjell, K., & Schwartz, H. A. (2022). Natural language analyzed with AI-based transformers predict traditional subjective well-being measures approaching the theoretical upper limits in accuracy. *Scientific Reports*, 12(1), 3918.

<https://doi.org/10.1038/s41598-022-07520-w>

Likert, R. (1932). A technique for the measurement of attitudes. *Arch. Psychol.*, 140(55).

Lovibond, P. F., & Lovibond, S. H. (1995). The structure of negative emotional states: Comparison of the Depression Anxiety Stress Scales (DASS) with the Beck Depression and Anxiety Inventories. *Behaviour Research and Therapy*, 33(3), 335–343.

[https://doi.org/10.1016/0005-7967\(94\)00075-u](https://doi.org/10.1016/0005-7967(94)00075-u)

- Lukat, J., Margraf, J., Lutz, R., Veld, W. M. van der, & Becker, E. S. (2016). Psychometric properties of the positive mental health scale (PMH-scale). *BMC Psychology*, 4(1), 8.  
<https://doi.org/10.1186/s40359-016-0111-x>
- Lutz, W., Schwartz, B., & Delgadillo, J. (2022). Measurement-Based and Data-Informed Psychological Therapy. *Annual Review of Clinical Psychology*, 18(1), 71–98.  
<https://doi.org/10.1146/annurev-clinpsy-071720-014821>
- Margraf, J., Cwik, J. C., Brachel, R. von, Suppiger, A., & Schneider, S. (2021). *DIPS open access 1.2: Diagnostisches interview bei psychischen störungen.*  
<https://doi.org/10.46586/rub.172.149>
- Matero, M., Idnani, A., Son, Y., Giorgi, S., Vu, H., Zamani, M., Limbachiya, P., Guntuku, S. C., & Schwartz, H. A. (2019). Suicide risk assessment with multi-level dual-context language and BERT. In K. Niederhoffer, K. Hollingshead, P. Resnik, R. Resnik, & K. Loveys (Eds.), *Proceedings of the sixth workshop on computational linguistics and clinical psychology* (pp. 39–44). Association for Computational Linguistics.  
<https://doi.org/10.18653/v1/W19-3005>
- Michalak, J., Kosfelder, J., Meyer, F., & Schulte, D. (2003). Messung des Therapieerfolgs. *Zeitschrift für Klinische Psychologie und Psychotherapie*, 32(2), 94–103.  
<https://doi.org/10.1026/0084-5345.32.2.94>
- Mohammadi, E., Amini, H., & Kosseim, L. (2019). CLaC at CLPsych 2019: Fusion of neural features and predicted class probabilities for suicide risk assessment based on online posts. In K. Niederhoffer, K. Hollingshead, P. Resnik, R. Resnik, & K. Loveys (Eds.), *Proceedings of the sixth workshop on computational linguistics and clinical psychology* (pp. 34–38). Association for Computational Linguistics.  
<https://doi.org/10.18653/v1/W19-3004>
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2022). Robust Speech Recognition via Large-Scale Weak Supervision. *arXiv Preprint arXiv: 2212.04356*. <https://arxiv.org/abs/2212.04356>
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267–288.

<https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is All you Need. *Advances in Neural Information Processing Systems, 30*.

Wampold, B. E., & Imel, Z. E. (2015). *The great psychotherapy debate: The evidence for what makes psychotherapy work*. Routledge.

World Health Organization. (2017). *Depression and other common mental disorders: Global health estimates*.

Zirikly, A., Resnik, P., Uzuner, Ö., & Hollingshead, K. (2019). CLPsych 2019 shared task: Predicting the degree of suicide risk in Reddit posts. In K. Niederhoffer, K. Hollingshead, P. Resnik, R. Resnik, & K. Loveys (Eds.), *Proceedings of the sixth workshop on computational linguistics and clinical psychology* (pp. 24–33). Association for Computational Linguistics. <https://doi.org/10.18653/v1/W19-3003>

Time-point	DU-DI	DU-Prä	KZT1-DU4	KZT1-DUPost	KZT2-DUPost	LZT1-DUPost	LZT2-DUPost	Kat6
<b>Explanation</b>	Pre-therapy, 4th contact	Pre-therapy, 6th contact	4th therapy session	12th therapy session	24th therapy session	45th therapy session	60th therapy session	6 month after therapy
<b>Diagnosis</b>	X							
<b>Demographics</b>		X						
<b>BSI</b>	X		X	X	X	X	X	X
<b>BDI-II</b>		X	X	X	X	X	X	X
<b>DASS-42</b>		X	X	X	X	X	X	X
<b>PMH</b>								
<b>CTQ</b>		X						
<b>CGI-S</b>		X						
<b>CGI-I</b>								
<b>Glob-Pt</b>								
<b>GAS</b>			X	X	X	X	X	X

**Figure 1**

*The step-by-step process of patient narrative analysis, from preprocessing to data evaluation.*

