

# Ethical Thematic and Topic Modelling Analysis of Sleep Concerns in a Social Media-Derived Suicidality Dataset

Martin Orr and Kirsten van Kessel

Auckland University of Technology, New Zealand

[martinorr521@gmail.com](mailto:martinorr521@gmail.com)

[kirsten.vankessel@aut.ac.nz](mailto:kirsten.vankessel@aut.ac.nz)

David Parry

Murdoch University, Australia

[david.parry@murdoch.edu.au](mailto:david.parry@murdoch.edu.au)

## Abstract

**Objective:** A thematic and topic modelling analysis of sleep concerns in a social media-derived, privacy-preserving, suicidality dataset. This forms the basis for an exploration of sleep as a potential computational linguistic signal in suicide prevention. **Background:** Suicidal ideation is a limited signal for suicide. Developments in computational linguistics and mental health datasets afford an opportunity to investigate additional signals and to consider the broader clinical ethical design implications. **Methodology:** A clinician-led integration of reflexive thematic analysis, with machine learning topic modelling (BERTopic), and the purposeful sampling of the University of Maryland Suicidality Dataset. **Results:** Sleep as a place of 1) refuge and escape, 2) revitalisation for exhaustion, and 3) risk and vulnerability were generated as core themes in an initial thematic analysis of 546 posts. BERTopic analysing 21,876 sleep references in 16791 posts facilitated the production of 40 topics that were clinically interpretable, relevant, and thematically aligned to a level that exceeded original expectations. Power and consent, privacy and synthetic representative data, validity and stochastic variability of results, and a co-designed and governed, multi-signal formulation perspective, are highlighted as key research and clinical issues.

## 1 Introduction

This paper reports on the thematic and topic modelling analysis of sleep concerns in a social media-derived, privacy-preserving, suicidality dataset. Key objectives are an exploration of 1) the role of sleep as a potential linguistic signal in suicide prevention formulation, and 2) the ethical and design opportunities and challenges, artificial intelligence (AI) and sensitive mental health datasets may afford to the global mental health community (Resnik et al., 2021; Shing et al., 2018, 2020; Zirikly et al.,

2019; Orr et al., 2022, 2023). This work arises from an academic program, centered on clinical best practice, leadership, and change and is the third in a series of papers focused on AI, ethics and suicide prevention (Orr et al., 2022, 2023).

There is increasing interest in the application of computational linguistics in suicide prevention. Suicide prevention threat detection, and guardian angel-type technology, are already deployed in social media. This raises significant ethical design, clinical effectiveness, and governance issues (Barnett and Torous, 2019; Bernert et al., 2020; Burke et al., 2019; Floridi and Cows, 2019; Ophir et al., 2022; Orr et al., 2023). Suicidal ideation is a relatively weak signal for suicide, and there is a need to research additional signals (Deisenhammer et al., 2009; Galynker et al., 2017; Yaseen et al., 2019). Sleep disturbance is associated with increased suicidal behaviour to the degree it may represent a modifiable risk factor and signal to inform suicide risk management formulation and intervention planning (Bishop et al., 2020; Bradford et al., 2021; Fernandes et al., 2021; Hamilton et al., 2023; Kalmbach et al., 2022; Liu et al., 2020; Miller and McCall, 2023; Shepard et al., 2023).

Computational linguistics is a branch of AI, associated with natural language processing. It has come to increasing clinical importance because of the rapid advancements and growth in transformers and generative AI (Javaid et al., 2023; Resnik et al., 2021).

Closely entwined with the use of AI in mental health, is the ethical utilisation and governance of sensitive data, including data mining and the curation of datasets. The data utilisation approach in this research combines the qualitative Braun and Clarke approach to thematic analysis, with quantitative machine learning topic modelling, to study a large social media suicidality dataset (Blei, 2012; Blei et al., 2003; Grootendorst, 2022; Braun and Clarke, 2006, 2019; Fast et al., 2016). There is

emerging interest in the potential research benefits of combining qualitative reflexive thematic analysis with computational methods. This includes how topic modelling could assist with more rapid familiarization and coding of large social media-based data sets, and contribute to the qualitative, nuanced, contextual, interpretative, and reflexive theme creation process (Gauthier et al., 2022; Gauthier and Wallace, 2022).

Suicide is a complex devastating event where there can be a struggle to make sense, find meaning, and understand the related computational basis of the mind, lived experience, intent, and decision-making. Computational neuroscience and psychiatry seek to bring understanding to complex human behaviours to optimize care and prevention (Hauser et al., 2022; Nordin et al., 2022).

Computational linguistics can play a role in data collection, formulation, and the creation of digital psychotherapeutic interventions. Either individually or combined with other neurophysiological, behavioural, and imaging techniques, computational linguistics may provide modelling insights into the neural basis of language (Bourguignon, 2022). This may contribute to the development, analysis, and detection of language biomarkers or signals of various states of cognitive and emotional processing, mental disorder, and behavioural risk.

Suicidal ideation is an important risk signal for potential psychological distress and completed suicide that requires timely assessment and intervention. However, although important, it is also a weakly predictive, frequently late, or non-presenting and unreliable signal (Deisenhammer et al., 2009; Galynker et al., 2017; Yaseen et al., 2019). Although there is only a limited correlation between suicidal ideation and suicide, those who have expressed suicidal ideation may be a significant target population to look for other signals operating with different contextual factors, frequency, prevalence, and time scales that may assist with the timely formulation of risk management and suicide prevention (Orr et al., 2022, 2023).

### **1.1 Sleep as a potential signal and intervention priority in suicide prevention**

Sleep is central to physical and mental health (Scott et al., 2021; Hertenstein et al., 2022; Harvey, 2022). There is significant evidence of sleep disturbance being correlated with suicidal behavior to the degree that it can be considered a potentially modifiable risk factor and predictive or prioritizing

signal for suicide prevention (Bishop et al., 2016, 2020; Blake and Allen, 2020; Bradford et al., 2021; Chaïb et al., 2020; Fernandes et al., 2021; Geofroy et al., 2021; Hamilton et al., 2023; Kalmbach et al., 2022; Kearns et al., 2020; Liu et al., 2020; McCall et al., 2019; McCall, 2022; Miller and McCall, 2023; Perlis et al., 2016; Pigeon et al., 2019; Shepard et al., 2023; Trockel et al., 2015; Tubbs et al., 2019).

Sleep disturbances from insomnia to nightmares to sleep-disordered breathing are associated with an increased risk of suicidal behavior (Porrás-Segovia et al., 2019; Prguda et al., 2023; Tubbs et al., 2019). However, there has been limited research targeting sleep as a suicide intervention and limited evidence of resultant benefit. Studies to date have typically targeted suicidal ideation with cognitive behavioural therapy for insomnia (CBTI) or hypnotic drugs and demonstrated enough selected evidence for improvement to warrant further research and consideration in treatment protocols (McCall et al., 2019; Pigeon et al., 2019; Trockel et al., 2015). Targeting suicidal ideation in research may be a useful proxy variable for completed suicide as there is a significant correlation between the two (Nock et al., 2008a,b). However, it is important to understand the nature and limitations of that correlation. Relatively few of those who voice suicidal ideation will go on to complete suicide, and many of those who complete suicide will not be known to have voiced or experienced suicidal ideation. If suicidal ideation is the only focus or signal sought, significant timely suicide prevention opportunities may be lost.

Night-time is a high-risk period for suicide. It is unclear to what degree this is due to a range of factors including circadian or sleep deprivation-related decrease in frontal lobe function or decreased serotonin levels, or due to loneliness or lack of support. Being alone and awake at night may decrease the potential for intervention or distraction from a distressed consciousness. Being awake during a period of brain hypofrontality may be associated with affective and cognitive dysfunction leading to emotionality, impulsivity, and impaired decision-making. Adolescents who are likely to be a key demographic utilising social media may be particularly vulnerable to sleep-related suicidal behaviour (Porrás-Segovia et al., 2019; Tubbs et al., 2019).

Each stage of sleep may have a key evolutionary function in helping process the cognitions, emotions, and actions of the day, clearing out waste,

and optimising mind and body. Lack of sleep can impair our ability to creatively problem solve, exercise impulse control, process and experience emotions, calm the mind, be empathetic, or see aesthetic beauty and positivity in the world (Peretti et al., 2019).

Sleep is a universal dynamic periodic function, associated with a range of neurophysiological, emotional, behavioral, and contextual parameters. Sleep can be measured multimodally, and the multiple signals integrated to get a greater understanding of sleep quantity, structure, and quality, and their relationships to context. In this work, the focus is on the computational linguistic analysis of the sleep signal that may be contained in social media text. A signal typically carries indicative information that may require or precipitate a specific response if when interpreted (either individually or combined with other signals or factors) it exhibits certain characteristics or meets a certain threshold. An individual's cognitions, emotional response, behaviours, and context in relation to sleep may be derived from how they write about it; and what and how much they write about sleep may be influenced by their mental and emotional state at that time and their personal, group, and cultural interpretation of the meaning and function of sleep and sleep disturbance.

Various forms of thematic analysis and phenomenological inquiry into the experience and understanding of sleep disturbance (insomnia, nightmares) and mental illness and suicide have previously been performed (Hochard et al., 2019; Klingaman et al., 2019; Littlewood et al., 2016; Luhaäär and Sisask, 2018). These are typically small in sample size, interview-based, and reflective in nature in the context of established mental disorder or suicidal behaviour.

If we can better linguistically define, characterise, or categorise the sleep experience and signal by and for both human and machine processes, it may contribute to the design of future suicide prevention research and interventions. A sleep signal may have utility in combination with other signals in terms of alerting, triage, clinical formulation, and treatment planning.

## 1.2 Identification of the University of Maryland Suicidality Dataset

There is an increasing focus on the application of data mining in combination with AI to enhance mental health service design and suicide prevention

(Berrouiguet et al., 2019; Lopez-Castroman et al., 2020; Schuerkamp et al., 2023; Wang, 2023).

An exploratory goal of the body of work to which this study relates is creating a conceptual linguistic sleep signal model that could contribute to the development of real-time, natural language processing empowered, social media and formulation-based suicide prevention. To align with this goal the modelling data utilised should emulate as much as feasible the timely naturalistic expression of the lived experience and context under study (Neubauer et al., 2019; Van Manen, 2017). Capturing live social media data in a naturalistic contextualised form from actively suicidal individuals, affords major ethical, clinical, and medicolegal issues, including issues of power and consent, and responsibility. Consultation with ethics and AI specialists concluded that live data was not ethically justifiable or essential for this exploratory research and that a pre-existing social media dataset should be identified. This led to the identification of the University of Maryland Suicidality Dataset. The dataset comprises the 11,129 users who between 2006 to August 31st, 2015, had posted on the subreddit r/suicidewatch and had posted 10 or more times in total across all of Reddit. Included are user posts, post ID, anonymised user ID, timestamp, subreddit, de-identified post title, and body. The dataset also has an equal number of controls who had not posted in r/suicidewatch. The r/suicidewatch subreddit focuses on individuals posting about their suicidal ideation and plans, and other users offering support. The total dataset has approximately 2 million documents. The dataset has expert-labelled, crowdsourced-labelled, and unlabelled subsets (Shing et al., 2018; Zirikly et al., 2019).

## 1.3 Reflexive Thematic Analysis

Thematic analysis is recognised as a reflexive interpretive approach to pattern or theme development across a dataset and for its utility in phenomenological or experiential qualitative research where there is a focus on the understanding of experience, meaning, and sensemaking (Braun and Clarke, 2019, 2021a,b).

Reflexive thematic analysis emphasizes the interaction of the researcher with the data in the qualitative creation of themes. Themes are inductively woven from codes, with fewer themes potentially illustrating more intricacy of the analytic thought. Though grounded in data, there is a need for reflex-

ive critical reflection, self-awareness, and recognition of prior knowledge, assumptions, and theories that may deductively influence the researcher's inductive sense-making process (Braun and Clarke, 2006; Clarke and Braun, 2018; Braun and Clarke, 2019).

The knowledge and experience of the primary researcher as a psychiatrist was utilised in the reflexive interpretation of data in this research (Braun and Clarke, 2021b; Ho et al., 2017; Neubauer et al., 2019; Pérez Vargas et al., 2020; Tomkins and Eatough, 2018; Van der Walt, 2020; Van Manen, 2017).

#### 1.4 Topic Modelling and BERTopic

Topic modelling involves a range of algorithmic techniques that are essentially quantitative in that the technology has no inherent sentient understanding of the text but can bring varying levels of prior knowledge and types of process to draw mathematical connections and create clusters of key characteristic terms that appear to be linked, and provide examples of the documents that best exemplify the links. The underlying quantitative paradigm is typically viewed as being one of discovery of latent topics or themes in the text across a body of documents. However typically and traditionally human interpretation is required to create sense and meaning around how the terms or words may be linked and be of utility and to appropriately guide or supervise the topic labelling (Al Moubayed et al., 2020; Blei et al., 2003; Blei, 2012; Chang et al., 2009; Kherwa and Bansal, 2019; Resnik et al., 2015; Vaswani et al., 2017).

This analysis utilises BERTopic a topic modelling tool that uses BERT (Bidirectional Encoder Representations from Transformers) embeddings (Devlin et al., 2018; Grootendorst, 2022).

#### 1.5 Central Question

The central question for the integrated reflexive thematic analysis and topic modelling was: what are the themes, topics, and key representative terms that may communicate or signal the experience, relationships, and meaning of sleep and sleep disturbance in those who have expressed suicidal ideation in social media text?

#### 1.6 Methodology

The study uses a mixed methods design (Johnson, 2017). The research involves the integration of

reflexive thematic analysis, BERTopic topic modelling, and the purposeful sampling of the University of Maryland Reddit Suicidality Dataset.

The Braun and Clarke reflexive thematic analysis process has six phases: 1) familiarisation with the data; 2) coding; 3) generating initial themes; 4) reviewing themes; 5) defining and naming themes; 6) writing up. The Braun and Clarke reflexive form of thematic analysis recognises how we bring our past knowledge, experience, and biases to the process of constructing patterns of meaning (Braun and Clarke, 2006, 2019, 2021a,b).

BERTopic creates topic representations in 3 major stages: 1) use of a pre-trained transformer language model to convert each document to its embedding mathematical representation; 2) optimisation of the embedding clustering process via reduction of the dimensionality; 3) topic representations are generated from the document clusters with a class-based variation of Term Frequency-Inverse Document Frequency (c-TF-IDF). Topic representations take the form of lists of keywords or terms that are most important, relevant, and characteristic of the topic, and BERTopic also creates a short collection of the most representative documents for each topic (Grootendorst, 2022).

## 2 Results

### 2.1 Thematic Analysis

The reflexive thematic analysis involved 546 posts that included the term sleep from 154 individuals in the expert risk-rated subset of 245 users. The thematic analysis utilised the posts (across the whole of Reddit) from the 245 expert risk-level-rated r/suicidewatch users. The posts were keyword searched looking for references to sleep. The themes generated in relation to the research question were sleep as a place of 1) refuge and escape, 2) risk and vulnerability, and 3) revitalization for exhaustion. These themes related to 1) seeking refuge and escape via sleep from the living nightmare and trauma of consciousness and physical and psychological pain; 2) feeling at risk and vulnerable to trauma and nightmares and sleep paralysis if enter sleep, or vulnerable in terms of where sleeping or who sleeping with, or vulnerable to being woken up by others and pets or vulnerable to insomnia and related anxiety, loneliness, pain, negative thoughts, constant arousal and being on edge and unable to switch the mind off with fear of missing out; 3) feeling constantly exhausted physi-

cally and psychologically and overwhelmed, tired of anticipatory anxious worry about the future, and ruminating worry about the past, feeling burnt out and seeking the revitalization of sleep.

## 2.2 Topic Modelling

The data subset for the topic modelling stage was derived by searching for the word sleep in all posts by users who had posted on r/suicidewatch in the University of Maryland dataset. This resulted in identifying 16791 posts by 5751 unique users. Those posts were then broken down into sentences, and sentences were selected that contained the word sleep. This identified 21,876 sentences.

BERTopic was utilised using a Google Colab that was customized to allow options in a range of parameters, including file selection, number of topic words, topic reduction, and seed topics. BERTopic provides a range of outputs, including key representative topic terms and documents, and visual representations of how the topics cluster and relate to each other and could potentially be merged or reduced. The BERTopic outputs in the Appendix section, Figure 1. (topic word scores), Table 2. (topic frequency count), and Figure 2. (hierarchical clustering) and the topic representative documents, were utilised by the psychiatrist first author for the topic labelling process (Table 1.). The thematic analysis was also drawn upon for a deeper interpretive and integrative understanding. The BERTopic outputs reported are with parameters set to 9 topic words, topic reduction to 40, and no guiding seeding. The non-seeded topic modelling process resulted in surprisingly clinically interpretable and relevant results, that generally supported and aligned with the core thematic analysis concepts. There are a range of qualitative and quantitative computational techniques by which BERTopic's outputs can be evaluated for topic coherence and diversity (Grootendorst, 2022). In this exploratory study the focus was on human domain expert interpretation.

The most relevant and important topic terms (with the highest topic c-TF-IDF scores), typically gave a readily interpretable guide to an appropriate, meaningful topic label e.g. Topic 10. Dreams and nightmares, Topic 13. Tired and exhausted and Topic 14. Pain, hurts, and sleep. These topics also aligned with the theme elements of exhaustion and vulnerability created in the thematic analysis. Sleep as escape was a major theme element in the thematic analysis. This was overtly captured in the

relatively low frequency Topic 36 where the most important and relevant term was escape and the representative documents directly refer to sleep being the only escape. However, escape was also an inherent, at least part, element of a range of other topics including the high-frequency topics 0. Want to sleep and 1. Sleeping pills to sleep. Though not immediately evident or definitive from the topic terms, reference to the representative documents indicated these topics capturing a want to sleep forever for some that aligned with the theme of escape and suicidality. That is references in the representative documents, to never wanting to wake up, or consuming a large supply of sleeping medication at once, were important potential suicidality signals. While in topics 0. and 1. the "forever" aspect only related to some, in Topic 24 it was the predominant term and feature.

The representative documents utilized in the topic labelling process, are not included in the Appendix because of the ethical requirement not to share verbatim potentially identifiable quotes and for data to be reported at the summative, coding, topic, and thematic level.

## 3 Discussion

This was exploratory research, with a focus on reporting and contextualising the analysis process and results, to a clinical /non-data scientist audience, to particularly highlight conceptual, ethical, and design issues for future research, and development. The following section aims to explore further a number of these issues.

### 3.1 Language and the Psyche

When considering the application of computational linguistics to mental health, it is important to also conceptually consider how language may relate to the functions of the mind. Language symbolically captures an individual's experience and conceptual interpretation of their world. Experience and interpretation are influenced by cultural and group norms. Language is integrally woven with brain function, and although the exact nature of the weave may be contentious, it may provide important data modelling insights to the understanding of the psyche, including rich multidimensional temporal and contextual parameters that may be difficult to access via other signals or means (Kompa, 2023; Li, 2022).

Topic	
-1: Miscellaneous (Outliers)	19: Gaming and Sleep
0: Want to sleep (forever)	20: Physical sensations and activity and sleep
1: Sleeping Pills to sleep (forever)	21: Sleep Apnea and Breathing and Sleep
2: Distressed and crying self to sleep	22: Texting and communication and sleep
3: Sleep period and schedule	23: Wish for more sleep
4: Name substitution and sleeping context	24: Want to sleep forever/permanently
5: Technology sleep mode	25: Mood disorder and sleep
6: Smoking, cannabis, drinking alcohol, and sleep	26: Prayer, meditation, and sleep
7: Eating and Sleep	27: Memory, focus, concentration and sleep
8: School and Sleep	28: Happiness and sleep
9: Sleeping outside home	29: Alarm clock and sleep
10: Dreams and nightmares	30: Motivation and sleep
11: Pets and sleep	31: Heat and temperature and sleep
12: R/nosleep stories	32: Cutting, self-harm and sleep
13: Tired and exhausted	33: Cuddling and sex and sleep
14: Pain, hurts and sleep	34: Sleep deprivation and health consequences
15: Sleeping locations	35: Scratching, itch and sleep
16: Mind and sleep	36: Sleep as escape
17: Sleep paralysis and neurological experiences	37: Nursing/breastfeeding and sleep
18: Music, Noise and Sleep	38: Hallucinations/psychosis and sleep

Table 1: Sleep Topic Labelling.

### 3.2 Ambivalence

Legally a declaration of suicide typically requires evidence of an intentional and knowing act. However, suicide may be characterized by ambivalence, conflicting cognitions, emotions, and behaviours, and a temporal perceived need to escape an overwhelmed or pained consciousness or sense of entrapment, rather than a specific knowing, reasoned intent to die (Orr et al., 2023). This presents challenges but also opportunities when considering a computational linguistic signal of the mind in relationship to suicidality, in terms of detection, formulation, guidance, and amplification in the direction of seeking help. Ambivalence is prevalent in the University of Maryland Suicidality Dataset and indeed could be considered inherent, as those posting are typically seeking some form of help, advice, and input from others. This could also be considered a limitation or caveat for the dataset in that the active group are actively signalling their risk, expressing suicidal ideation, and reaching out. Those who don't express suicidal ideation and don't reach out may differ, adding weight to the case for additional signals and methods and channels for detection.

### 3.3 Sensitive data, stochasticity, and variability concerns

Two of the major concerns and limitations of the use of computational linguistics in clinical practice and research are, 1. concerns around the privacy, security, and governance of sensitive data and 2. the validity and variability and related safety of outputs. Validity, stability, and reproducibility are key concerns of topic model-based content analysis (Hoyle et al., 2022). The stochastic nature of computational linguistic processes may contribute to this variability (Javaid et al., 2023).

BERTopic is stochastic with variability in outputs on each run (mainly related to UMAP) (McInnes et al., 2018; Grootendorst, 2022). Stochasticity is a central feature of large language models and generative AI, which makes it an increasingly relevant concept for clinical research and related awareness. Stochasticity or probabilistic variability in potential outputs can augment and amplify human creativity, engagement, and brainstorming. However, in a clinical context the resultant variability in outputs (and related perceived confabulation or hallucinations with generative AI), can be perceived as unsafe.

Safe clinical and research utilization of pre-trained transformers and generative AI will require

significant development in education, prompt engineering training, stochastic temperature control, guard railing, and fine-tuning. This may enhance both the factual certainty and relevancy potential of the technology, while also being able to leverage the creative potential of stochasticity, to assist with resolving complex problems like suicide prevention.

### **3.4 Thematic analysis and topic modelling**

The labelled topics offer a clinically valid and relevant general range of subdomains for issues related to sleep and sleep disturbance including aligning with key codes and themes created during the thematic analysis process.

The option of seed topics was included in the BERTopic Colab design as there was a high expectation that significant guidance and interpretation would be required to get any form of meaningful results. In Guided BERTopic seeds are used to nudge towards the creation of particular topics, but if they do not exist within the dataset they will not be modelled. The original expectation was that seed topics influenced by the earlier thematic analysis work would be required for meaningful and aligned results.

As it transpired the topic modelling exceeded expectations in terms of clinically interpretable and meaningful outputs, even with minimal guidance. The reflexive thematic analysis process still served an important purpose, in enhancing a deeper understanding of the dataset and facilitating the interpretation and labelling of the outputs. The thematic outputs from reflexive thematic analysis, are personal subjective qualitative constructs of the data, as it relates to the research question. The construction of themes in reflexive thematic analysis, and computational linguistic processes such as topic modelling, have a range of similarities in approach, in terms of pattern recognition, baseline weights based on prior learning, contextualisation, and a capacity to iteratively fine-tune.

### **3.5 Linguistic signal formulation**

Signals in psychiatry can be characterized by complexity, noise, dissonance, probability, uncertainty, ambiguity, and ambivalence. This is similarly true for language signals that are subject to significant semantic and pragmatic interpretation complexity, and varying levels of contextual, and cultural modifiers and abstractive symbolism in deriving meaning, sentiment, and intent.

It is important to take a formulation approach to signals that recognises the importance of complexity, context, and culture and the need to dynamically consider and weigh all other factors or signals. Humans signal their thoughts, sentiments, and intent in a range of complex neurophysiological, behavioural, and natural language ways. Formulation is core to clinical mental health practice and has a factor weighing, pattern recognition, and modelling focus. Formulation recognises that the explanation for human behaviour can be complex, contextual, and contingent (Orr et al., 2022).

### **3.6 AI, data and the ethics of research and development and power and consent**

There is an increasing focus on the need for an ethical overview of social media and AI research. There is a move away from considering all public data as exempt from ethical board oversight and more focus on the complexities of consent, defining private and public, anonymity, sensitive data and vulnerable populations, and minimising bias and algorithmic harm (Benton et al., 2017; British Psychological Society, 2017; Chiauzzi and Wicks, 2019; NEAC, 2019; Pagoto and Nebeker, 2019; Townsend and Wallace, 2016; Organization, 2023).

The topic representations created by BERTopic can be fine-tuned and labelled by a range of methods including via Open AI's API to ChatGPT (Grootendorst, 2022). To respect and protect the principle and ethics requirement of not sharing real data and quotes from the dataset, this process was not utilized. Future challenges include developing methods that respect and protect the sensitive nature of the data while leveraging and evolving the benefits and interpretive and generative nuance that large language models may afford.

The AI algorithmic and global mental health realms have both been subject to debate and scrutiny around ethics, power dynamics, bias, agency, consent and control. This includes concerns around AI and global mental health running the risk of representing new forms of colonialism and imperialism. Those who control data and algorithms may have undue power and influence and knowingly or unknowingly not act in the best ethical interests of an individual or community. (Beresford and Rose, 2023; Birhane, 2023; Pendse et al., 2022). AI may be associated with dual-use and algorithmic harm including unintentional iatrogenic harm if used for clinical purposes and not appropriately researched, designed or governed for different

contexts, cultures, and customised needs.

### 3.7 Quotes and Sensitive Data

Qualitative research has traditionally used, and indeed required verbatim quotes, as a form of evidence of quality and rigour. Research using social media data sets will increasingly challenge that tradition, in that verbatim quotes can be frequently identified with a simple search engine and triangulation approach. There is a broad range of terms that may be utilised for indicative quotes, that do not contain the original words or syntax, but seek to convey the original content, sentiment, and intent. These terms include clustered, blended, aggregated, combined, composite, collective, deidentified, spun, paraphrased, bundled, amalgamated, illustrative, characteristic, indicative, representative and synthetic, and synthesized (Hemphill et al., 2022; Kasal et al., 2023; Proferes et al., 2021; Reagle, 2022; Winter and Gundur, 2022; Zimmer, 2020).

There will be an increasing range of technological and generative AI options to carry out this deidentification function with various degrees of parsing, paraphrasing, and production. However indicative non-verbatim quotes may not capture the emotion, pain, beauty, poignancy, pragmatic metaphorical abstraction, and personal poetry of the original where an individual has crafted their personal experience into words that they want to cathartically share with others. Similarly, the more data that is produced or created by generative AI, the more this may decrease the capacity to validly understand the authentic expression of human experience by drowning out and indeed shaping the expression of that experience.

Table 3. in the Appendix provides examples of composite synthetic quotes, that aim to illustrate the type of content that was the basis of the reflexive thematic analysis generated themes, without infringing the ethical undertaking not to use verbatim quotes.

## 4 Conclusion

This paper reported on the thematic and computational linguistic topic modelling analysis of sleep concerns in the University of Maryland Suicidality Dataset. This was multidisciplinary, exploratory, foundational work, that had the broader aim of highlighting some of the conceptual, ethical, and design opportunities and challenges artificial in-

telligence and mental health datasets may afford. The reflexive thematic analysis produced three core themes; sleep as a place of refuge and escape, risk and vulnerability, and revitalization for exhaustion. BERTopic was utilized to produce 40 topics with representative key terms and documents. The combined thematic analysis, and topic modelling process, resulted in clinically interpretable, relevant, and aligned results that exceeded initial expectations.

This is the third in a series of papers focused on AI and suicide prevention. Central series themes have been the complexities and contentions of suicide prediction, the related central role of formulation in clinical practice, and how the computational linguistic detection, development, and integration of relevant signals may contribute to enhancing the formulation and intervention planning process.

Sleep is a potentially useful linguistic signal in AI-based suicide risk formulation and intervention planning. Establishing sleep themes, topics, and key terms represents an initial exploratory development stage. A deployed AI-enhanced social media-based system that could detect, and utilize a linguistic sleep signal would need significant ethical co-design and governance. Research and development would require an iterative multistage, multimodal, multisignal contextually and culturally aware integrated formulation approach. Sleep may be considered as both a signal and an intervention and more conceptually as a preferable escape for a traumatised, overwhelmed, and exhausted consciousness. Sleep may help with cognitive and emotional processing, decreasing impulsivity, and increasing the capability to see a positive path through. An advantage of focusing on sleep as a key signal is that it can be both a transdiagnostic indicator of illness and vulnerability, and a transdiagnostic positive intervention, maximising the opportunities for benefit and optimisation of scarce resource utilisation.

Suicide is a complex, multifactorial low-base rate event, where there are significant risks and limitations in prediction and particularly attaching specific predictive or intervention power to any standalone factor. Sleep and sleep disturbance may have a significant role as a sensitive indicator of human distress and arousal and mental health vulnerability and as a signal that further assessment and intervention are required. Sleep and sleep disturbance signal data may be coded, labelled, weighed, and used in the formulation both as an indicator of transdiagnostic emotional distress and arousal,



and to contribute towards specific mental disorder and sleep disorder interventions. It may also be more acceptable or in keeping with cultural, group, or personal experiential norms to talk about sleep disturbance and related exhaustion than to report depression or suicidal ideation.

Data relating to suicide is highly sensitive, and privacy-preserving datasets may afford clinical safety, medicolegal, and ethical benefits. Social media data mining affords an opportunity, to improve the computational understanding of human behavior and provide insights into user mental models. This may enhance psychological formulation, targeted needs segmentation, and personalized timely, user experience, engagement, recommendation, and intervention planning. This is particularly important in suicide prevention. In terms of practical clinical contributions, arising from this research, the findings have contributed to the psychiatrist first author's design of digital intake forms for a specialist sleep clinic. The future aim would be to develop a multimodal AI-enhanced assessment and formulation system, that was also capable of assisting with the therapeutic revision of memory and narrative associated with traumatic nightmares. Nightmares are associated with suicidality and were a significant finding of both the thematic analysis and topic modelling aspects of this research.

In terms of methodological contributions, the integration of reflexive thematic analysis and machine learning topic modelling could benefit other researchers working with large datasets that require scale, speed, and interpretive nuance in analysis. The integrated approach may contribute to the explicability or understanding of the algorithmic process and thematic and topic results.

As the AI field moves from Large Language Models to Merged, Multimodal Models (Triple M's), computational linguistics should retain a central role in sensing, shaping, augmenting, and amplifying the human psyche toward creative effective action and outputs. There is a need to research how large language models and generative AI could be used in fine-tuning, and clinically relevant analysis and interventions, in a way that addresses the sensitive nature of potential mental health and risk-related data and the associated medicolegal, ethical, clinical effectiveness, and safety issues.

The stated vision for the CLPsych community is to improve interdisciplinary knowledge exchange, foster collaboration, and increase the visibility of mental health as a problem domain in natural lan-

guage processing. Similarly important is improving the visibility and accessibility of computational linguistics as an opportunity domain in clinical practice. The rapid rise in pre-trained transformers and generative AI has increased the need for clinician knowledge, education, and engagement in computational linguistics. Also highlighted by these developments is a broader need for multi-stakeholder research, co-design and representative governance to enhance the capacity and propensity for benefit optimization and harm minimisation. Key related areas for future research and development highlighted by this paper include building a shared understanding and approach to, risk-benefit analysis, power and consent, formulation, stochastics, and sensitive and representative data.

By focusing on one signal (sleep), one dataset, and one technique (topic modelling) a core aim of this research is a greater shared conceptual understanding of the opportunities and challenges presented by computational linguistics to the global mental health community. Global mental health may be particularly suited and vulnerable to developments in computational linguistics due to the central role of language in the psyche and culture and the need for contextual weighting in making sense of and shaping complex human experiences.

The computational linguistic and global mental health communities need to engage with these challenges and opportunities with a high degree of shared ethical, conceptual, contextual, cultural, and power dynamic awareness. This paper has aimed to enhance that awareness and contribute to the developmental capacity of the CLPsych community to pursue its mission to reduce emotional suffering and suicide.

## Acknowledgements

We would like to acknowledge the assistance of Prof. Philip Resnik and the American Association of Suicidology in making the "University of Maryland Reddit Suicidality Dataset, Version 2" available.

**Ethics approval:** Locality academic ethics committee approval was granted on 21st Sept. 2020. University of Maryland/American Association of Suicidology approval for dataset usage granted 29th Sept. 2020. Conflicts and funding: Nil  
**The dual approval requirements are described at:** [https://users.umiacs.umd.edu/~resnik/umd\\_reddit\\_suicidality\\_dataset.html](https://users.umiacs.umd.edu/~resnik/umd_reddit_suicidality_dataset.html)

**Limitations statement:** The relationships between sleep, mental health, and suicide are complex and contentious. As highlighted in the paper this was exploratory research, that formed part of a body of work with a focus on ethical clinical design, and the opportunities and challenges of computational linguistics contributing to enhancing mental health care and research, specifically suicide prevention. There was a concern from the outset that limitations of the dataset and BERTopic pre-trained embeddings and processes could lead to few interpretable, clinically recognizable, or useful topics. The results exceeded expectations in terms of potential clinical utility and interest. The topic modelling results reported had a focus on illustrating potential clinical utility with a view to clinician engagement. It is well recognized that they represent just an initial analytic and developmental step in the ongoing evolution of safe and effective risk and clinical formulation and intervention systems. Caveats in terms of ethics, dynamic user, temporal, situational, and language representation and understanding, misinformation, bias, contextual and cultural nuance, generalisability, reproducibility and validity, and lack of ground truth around risk and outcome, and cause and correlation remain.

## References

- Noura Al Moubayed, Stephen McGough, and Bashar Awwad Shiekh Hasan. 2020. Beyond the topics: How deep learning can improve the discriminability of probabilistic topic modelling. *PeerJ Computer Science*, 6:e252.
- Ian Barnett and John Torous. 2019. Ethics, transparency, and public health at the intersection of innovation and Facebook’s suicide prevention efforts.
- Adrian Benton, Glen Coppersmith, and Mark Dredze. 2017. Ethical research protocols for social media health research. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 94–102.
- Peter Beresford and Diana Rose. 2023. Decolonising global mental health: The role of Mad Studies. *Cambridge Prisms: Global Mental Health*, 10:e30.
- Rebecca A Bernert, Amanda M Hilberg, Ruth Melia, Jane Paik Kim, Nigam H Shah, and Freddy Abnoui. 2020. Artificial intelligence and suicide prevention: A systematic review of machine learning investigations. *International Journal of Environmental Research and Public Health*, 17(16):5929.
- Sofian Berrouguet, María Luisa Barrigón, Jorge Lopez Castroman, Philippe Courtet, Antonio Artés-Rodríguez, and Enrique Baca-García. 2019. Combining mobile-health (mHealth) and artificial intelligence (AI) methods to avoid suicide attempts: The Smartcrises study protocol. *BMC psychiatry*, 19(1):1–9.
- Abeba Birhane. 2023. *Algorithmic Colonization of Africa*. In *Imagining AI: How the World Sees Intelligent Machines*. Oxford University Press.
- Todd M Bishop, Kelsey V Simons, Deborah A King, and Wilfred R Pigeon. 2016. Sleep and suicide in older adults: An opportunity for intervention. *Clinical Therapeutics*, 38(11):2332–2339.
- Todd M Bishop, Patrick G Walsh, Lisham Ashrafioun, Jill E Lavigne, and Wilfred R Pigeon. 2020. Sleep, suicide behaviors, and the protective role of sleep medicine. *Sleep Medicine*, 66:264–270.
- Matthew J Blake and Nicholas B Allen. 2020. Prevention of internalizing disorders and suicide via adolescent sleep interventions. *Current Opinion in Psychology*, 34:37–42.
- David M Blei. 2012. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent Dirichlet Allocation. *The Journal of Machine Learning Research*, 3:993–1022.
- Nicolas J Bourguignon. 2022. The emergence of language in the human mind and brain — insights from the neurobiology of language, thought and action. *Psychological Review*.
- Daniel RR Bradford, Stephany M Biello, and Kirsten Russell. 2021. Insomnia symptoms mediate the association between eveningness and suicidal ideation, defeat, entrapment, and psychological distress in students. *Chronobiology International*, 38(10):1397–1408.
- Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology*, 3(2):77–101.
- Virginia Braun and Victoria Clarke. 2019. Reflecting on reflexive thematic analysis. *Qualitative Research in Sport, Exercise and Health*, 11(4):589–597.
- Virginia Braun and Victoria Clarke. 2021a. Can I use TA? Should I use TA? Should I not use TA? Comparing reflexive thematic analysis and other pattern-based qualitative analytic approaches. *Counselling and Psychotherapy Research*, 21(1):37–47.
- Virginia Braun and Victoria Clarke. 2021b. One size fits all? What counts as quality practice in (reflexive) thematic analysis? *Qualitative Research in Psychology*, 18(3):328–352.
- British Psychological Society. 2017. Ethics guidelines for internet-mediated research. *Leicester, UK: British Psychological Society*.

- Taylor A Burke, Brooke A Ammerman, and Ross Jacobucci. 2019. The use of machine learning in the study of suicidal and non-suicidal self-injurious thoughts and behaviors: A systematic review. *Journal of Affective Disorders*, 245:869–884.
- Laurent Stephane Chaïb, Alejandro Porras Segovia, Enrique Baca-García, and Jorge Lopez-Castroman. 2020. Ecological studies of sleep disturbances during suicidal crises. *Current Psychiatry Reports*, 22:1–8.
- Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-Graber, and David Blei. 2009. Reading tea leaves: How humans interpret topic models. *Advances in Neural Information Processing Systems*, 22.
- E. Chiauzzi and P. Wicks. 2019. [Digital trespass: Ethical and terms-of-use violations by researchers accessing data from an online patient community](#). *J Med Internet Res*, 21(2):e11985.
- Victoria Clarke and Virginia Braun. 2018. Using thematic analysis in counselling and psychotherapy research: A critical reflection. *Counselling and Psychotherapy Research*, 18(2):107–110.
- Eberhard A Deisenhammer, Chy-Meng Ing, Robert Strauss, Georg Kemmler, Hartmann Hinterhuber, and Elisabeth M Weiss. 2009. The duration of the suicidal process: How much time is left for intervention between consideration and accomplishment of a suicide attempt? *Journal of Clinical Psychiatry*, 70(1):19.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Ethan Fast, Binbin Chen, and Michael S Bernstein. 2016. Empath: Understanding topic signals in large-scale text. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 4647–4657.
- Sara N Fernandes, Emily Zuckerman, Regina Miranda, and Argelinda Baroni. 2021. When night falls fast: Sleep and suicidal behavior among adolescents and young adults. *Child and Adolescent Psychiatric Clinics*, 30(1):269–282.
- Luciano Floridi and Josh Cowls. 2019. [A unified framework of five principles for AI in society](#). *Harvard Data Science Review*.
- Igor Galynker, Zimri S Yaseen, Abigail Cohen, Ori Benhamou, Mariah Hawes, and Jessica Briggs. 2017. [Prediction of suicidal behavior in high risk psychiatric patients using an assessment of acute suicidal state: The suicide crisis inventory](#). *Depression and Anxiety*, 34(2):147–158.
- Robert P Gauthier, Mary Jean Costello, and James R Wallace. 2022. “I will not drink with you today”: A topic-guided thematic analysis of addiction recovery on Reddit. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–17.
- Robert P Gauthier and James R Wallace. 2022. The computational thematic analysis toolkit. *Proceedings of the ACM on Human-Computer Interaction*, 6(GROUP):1–15.
- Pierre A Geoffroy, Maria A Oquendo, Philippe Courtet, Carlos Blanco, Mark Olfson, Hugo Peyre, Michel Lejoyeux, Frederic Limosin, and Nicolas Hoertel. 2021. Sleep complaints are associated with increased suicide risk independently of psychiatric disorders: Results from a national 3-year prospective study. *Molecular Psychiatry*, 26(6):2126–2136.
- Maarten Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794*.
- Jessica L Hamilton, Aliona Tsypes, Jamie Zelazny, Craig JR Sewall, Noelle Rode, John Merranko, David A Brent, Tina R Goldstein, and Peter L Franzen. 2023. Sleep influences daily suicidal ideation through affective reactivity to interpersonal events among high-risk adolescents and young adults. *Journal of Child Psychology and Psychiatry*, 64(1):27–38.
- Allison G Harvey. 2022. Treating sleep and circadian problems to promote mental health: Perspectives on comorbidity, implementation science and behavior change. *Sleep*, 45(4):zsac026.
- Tobias U Hauser, Vasilisa Skvortsova, Munmun De Choudhury, and Nikolaos Koutsouleris. 2022. The promise of a model-based psychiatry: Building computational models of mental ill health. *The Lancet Digital Health*, 4(11):e816–e828.
- Libby Hemphill, Angela Schöpke-Gonzalez, and Anmol Panda. 2022. Comparative sensitivity of social media data and their acceptable use in research. *Scientific Data*, 9(1):643.
- Elisabeth Hertenstein, Ersilia Trinca, Marina Wunderlin, Carlotta L Schneider, Marc A Züst, Kristoffer D Fehér, Tanja Su, Annemieke v Straten, Thomas Berger, Chiara Baglioni, et al. 2022. Cognitive behavioral therapy for insomnia in patients with mental disorders and comorbid insomnia: A systematic review and meta-analysis. *Sleep Medicine Reviews*, 62:101597.
- Ken HM Ho, Vico CL Chiang, and Doris Leung. 2017. Hermeneutic phenomenological analysis: The ‘possibility’ beyond ‘actuality’ in thematic analysis. *Journal of Advanced Nursing*, 73(7):1757–1766.
- Kevin D Hochard, Sam Ashcroft, Janine Carroll, Nadja Heym, and Ellen Townsend. 2019. Exploring thematic nightmare content and associated self-harm risk. *Suicide and Life-Threatening Behavior*, 49(1):64–75.

- Alexander Miserlis Hoyle, Pranav Goel, Rupak Sarkar, and Philip Resnik. 2022. [Are neural topic models broken?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5321–5344, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mohd Javaid, Abid Haleem, and Ravi Pratap Singh. 2023. ChatGPT for healthcare services: An emerging stage for an innovative perspective. *BenchCouncil Transactions on Benchmarks, Standards and Evaluations*, 3(1):100105.
- R Burke Johnson. 2017. Dialectical pluralism: A meta-paradigm whose time has come. *Journal of Mixed Methods Research*, 11(2):156–173.
- David A Kalmbach, Philip Cheng, Brian K Ahmedani, Edward L Peterson, Anthony N Reffi, Chaewon Sagong, Grace M Seymour, Melissa K Ruprich, and Christopher L Drake. 2022. Cognitive-behavioral therapy for insomnia prevents and alleviates suicidal ideation: Insomnia remission is a suicidolytic mechanism. *Sleep*, 45(12):zsac251.
- Alexandr Kasal, Roksana Táborská, Laura Juríková, Alexander Grabenhofer-Eggerth, Michaela Pichler, Beate Gruber, Hana Tomášková, and Thomas Niederkrotenthaler. 2023. Facilitators and barriers to implementation of suicide prevention interventions: Scoping review. *Cambridge Prisms: Global Mental Health*, 10:e15.
- Jaclyn C Kearns, Daniel DL Coppersmith, Angela C Santee, Catherine Insel, Wilfred R Pigeon, and Catherine R Glenn. 2020. Sleep problems and suicide risk in youth: A systematic review, developmental framework, and implications for hospital treatment. *General Hospital Psychiatry*, 63:141–151.
- Pooja Kherwa and Poonam Bansal. 2019. Topic modeling: A comprehensive review. *EAI Endorsed Transactions on Scalable Information Systems*, 7(24).
- Elizabeth A Klingaman, Alicia Lucksted, Eric S Crosby, Yelena Blank, and Elana Schwartz. 2019. A phenomenological inquiry into the experience of sleep: Perspectives of US military veterans with insomnia and serious mental illness. *Journal of Sleep Research*, 28(4):e12833.
- Nikola A Kompa. 2023. Inner speech and ‘pure’ thought – do we think in language? *Review of Philosophy and Psychology*, pages 1–18.
- Jing Li. 2022. Relationship between language and thought: Linguistic determinism, independence, or interaction? *Journal of Contemporary Educational Research*, 6(5):32–37.
- Donna L Littlewood, Patricia Gooding, Simon D Kyle, Daniel Pratt, and Sarah Peters. 2016. Understanding the role of sleep in suicide risk: Qualitative interview study. *BMJ open*, 6(8):e012113.
- Richard T Liu, Alexandra H Bettis, and Taylor A Burke. 2020. [Characterizing the phenomenology of passive suicidal ideation: A systematic review and meta-analysis of its prevalence, psychiatric comorbidity, correlates, and comparisons with active suicidal ideation.](#) *Psychological Medicine*, 50(3):367–383.
- Jorge Lopez-Castroman, Bilel Moulahi, Jérôme Azé, Sandra Bringay, Julie Deninotti, Sebastien Guillaume, and Enrique Baca-Garcia. 2020. Mining social networks to improve suicide prevention: A scoping review. *Journal of Neuroscience Research*, 98(4):616–625.
- Kätlin Luhaäär and Merike Sisask. 2018. Pathways to attempted suicide as reflected in the narratives of people with lived experience. *Religions*, 9(4):137.
- William V McCall. 2022. Targeting insomnia symptoms as a path to reduction of suicide risk: The role of Cognitive Behavioral Therapy for insomnia (CBT-I).
- William V McCall, Ruth M Benca, Peter B Rosenquist, Nagy A Youssef, Laryssa McCloud, Jill C Newman, Doug Case, Meredith E Rumble, Steven T Szabo, Marjorie Phillips, et al. 2019. Reducing suicidal ideation through insomnia treatment (REST-IT): a randomized clinical trial. *American Journal of Psychiatry*, 176(11):957–965.
- Leland McInnes, John Healy, and James Melville. 2018. UMAP: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*.
- Brian J Miller and William V McCall. 2023. Meta-analysis of insomnia, suicide, and psychopathology in schizophrenia. *Current Opinion in Psychiatry*, 36(3):156–165.
- NEAC. 2019. National ethical standards for health and disability research and quality improvement. *Wellington: Ministry of Health*.
- Brian E Neubauer, Catherine T Witkop, and Lara Varpio. 2019. How phenomenology can help us learn from the experiences of others. *Perspectives on Medical Education*, 8:90–97.
- Matthew K Nock, Guilherme Borges, Evelyn J Bromet, Jordi Alonso, Matthias Angermeyer, Annette Beautrais, Ronny Bruffaerts, Wai Tat Chiu, Giovanni De Girolamo, Semyon Gluzman, et al. 2008a. Cross-national prevalence and risk factors for suicidal ideation, plans and attempts. *The British journal of psychiatry*, 192(2):98–105.
- Matthew K Nock, Guilherme Borges, Evelyn J Bromet, Christine B Cha, Ronald C Kessler, and Sing Lee. 2008b. Suicide and suicidal behavior. *Epidemiologic Reviews*, 30(1):133.
- Noratikah Nordin, Zurinahni Zainol, Mohd Halim Mohd Noor, and Lai Fong Chan. 2022. Suicidal behaviour prediction models using machine learning

- techniques: A systematic review. *Artificial Intelligence in Medicine*, page 102395.
- Yaakov Ophir, Refael Tikochinski, Anat Brunstein Klomek, and Roi Reichart. 2022. The hitchhiker's guide to computational linguistics in suicide prevention. *Clinical Psychological Science*, 10(2):212–235.
- World Health Organization. 2023. *Regulatory considerations on artificial intelligence for health*. World Health Organization.
- Martin Orr, Kirsten Van Kessel, and David Parry. 2022. The ethical role of computational linguistics in digital psychological formulation and suicide prevention. In *Proceedings of the Eighth Workshop on Computational Linguistics and Clinical Psychology*.
- Martin Orr, Kirsten Van Kessel, and David Parry. 2023. Ethical suicide prevention in an artificial intelligence driven society. *Journal of Ethics in Mental Health*, 11.
- Sherry Pagoto and Camille Nebeker. 2019. [How scientists can take the lead in establishing ethical practices for social media research](#). *Journal of the American Medical Informatics Association*, 26(4):311–313.
- Sachin R Pendse, Daniel Nkemelu, Nicola J Bidwell, Sushrut Jadhav, Soumitra Pathare, Munmun De Choudhury, and Neha Kumar. 2022. From treatment to healing: Envisioning a decolonial digital mental health. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–23.
- Sara Peretti, Daniela Tempesta, Valentina Socci, Maria C Pino, Monica Mazza, Marco Valenti, Luigi De Gennaro, Cinzia Di Dio, Antonella Marchetti, and Michele Ferrara. 2019. The role of sleep in aesthetic perception and empathy: A mediation analysis. *Journal of Sleep Research*, 28(3):e12664.
- John Jairo Pérez Vargas, Johan Andrés Nieto Bravo, and Juan Esteban Santamaría Rodríguez. 2020. Hermeneutics and phenomenology in human and social sciences research. *Civilizar Ciencias Sociales y Humanas*, 20(38):137–144.
- Michael L Perlis, Michael A Grandner, Subhajt Chakravorty, Rebecca A Bernert, Gregory K Brown, and Michael E Thase. 2016. Suicide and sleep: Is it a bad thing to be awake when reason sleeps? *Sleep Medicine Reviews*, 29:101–107.
- Wilfred R Pigeon, Jennifer S Funderburk, Wendi Cross, Todd M Bishop, and Hugh F Crean. 2019. Brief CBT for insomnia delivered in primary care to patients endorsing suicidal ideation: A proof-of-concept randomized clinical trial. *Translational Behavioral Medicine*, 9(6):1169–1177.
- Alejandro Porrás-Segovia, María M Pérez-Rodríguez, Pilar López-Esteban, Philippe Courtet, Jorge López-Castromán, Jorge A Cervilla, Enrique Baca-García, et al. 2019. Contribution of sleep deprivation to suicidal behaviour: A systematic review. *Sleep Medicine Reviews*, 44:37–47.
- Emina Prguda, Justine Evans, Sarah McLeay, Madeline Romaniuk, Andrea J Phelps, Kerri Lewis, Kelly Brown, Gina Fisher, Fraser Lowrie, Elise Saunders-Dow, et al. 2023. Posttraumatic sleep disturbances in veterans: A pilot randomized controlled trial of cognitive behavioral therapy for insomnia and imagery rehearsal therapy. *Journal of Clinical Psychology*, 79(11):2493–2514.
- Nicholas Proferes, Naiyan Jones, Sarah Gilbert, Casey Fiesler, and Michael Zimmer. 2021. [Studying Reddit: A systematic overview of disciplines, approaches, methods, and ethics](#). *Social Media + Society*, 7(2):20563051211019004.
- Joseph Reagle. 2022. Disguising Reddit sources and the efficacy of ethical research. *Ethics and Information Technology*, 24(3):41.
- Philip Resnik, William Armstrong, Leonardo Claudino, Thang Nguyen, Viet-An Nguyen, and Jordan Boyd-Graber. 2015. Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 99–107.
- Philip Resnik, April Foreman, Michelle Kuchuk, Katherine Musacchio Schafer, and Beau Pinkham. 2021. Naturally occurring language as a source of evidence in suicide prevention. *Suicide and Life-Threatening Behavior*, 51(1):88–96.
- Ryan Schuerkamp, Luke Liang, Ketra L Rice, and Philippe J Giabbanelli. 2023. Simulation models for suicide prevention: A survey of the state-of-the-art. *Computers*, 12(7):132.
- Alexander J Scott, Thomas L Webb, Marrison Martyn-St James, Georgina Rowse, and Scott Weich. 2021. Improving sleep quality leads to better mental health: A meta-analysis of randomised controlled trials. *Sleep Medicine Reviews*, 60:101556.
- Christopher A Shepard, Katrina A Rufino, Jaehoon Lee, Tiffany Tran, Kieran Paddock, Chester Wu, John M Oldham, Sanjay J Mathew, and Michelle A Patriquin. 2023. Nighttime sleep quality and daytime sleepiness predicts suicide risk in adults admitted to an inpatient psychiatric hospital. *Behavioral Sleep Medicine*, 21(2):129–141.
- Han-Chin Shing, Suraj Nair, Ayah Zirikly, Meir Friedenberg, Hal Daumé III, and Philip Resnik. 2018. Expert, crowdsourced, and machine assessment of suicide risk via online postings. In *Proceedings of the Fifth Workshop on Computational Linguistics and Clinical Psychology: From Keyboard to Clinic*, pages 25–36.

- Han-Chin Shing, Philip Resnik, and Douglas W Oard. 2020. A prioritization model for suicidality risk assessment. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8124–8137.
- Leah Tomkins and Virginia Eatough. 2018. Hermeneutics: Interpretation, understanding and sense-making. *SAGE Handbook of Qualitative Business and Management Research Methods*, pages 185–200.
- Leanne Townsend and Claire Wallace. 2016. Social media research: A guide to ethics. *University of Aberdeen*, pages 1–16.
- Mickey Trockel, Bradley E Karlin, C Barr Taylor, Gregory K Brown, and Rachel Manber. 2015. Effects of cognitive behavioral therapy for insomnia on suicidal ideation in veterans. *Sleep*, 38(2):259–265.
- Andrew S Tubbs, Michael L Perlis, and Michael A Grandner. 2019. Surviving the long night: The potential of sleep health for suicide prevention. *Sleep Medicine Reviews*, 44:83.
- Johannes L Van der Walt. 2020. Interpretivism-Constructivism as a research method in the humanities and social sciences – More to it than meets the eye. *International Journal of Philosophy and Theology*, 8(1):59–68.
- Max Van Manen. 2017. Phenomenology in its original sense. *Qualitative Health Research*, 27(6):810–825.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*, 30.
- Zizhan Wang. 2023. Review on data mining and data analysis method for adolescent suicide problem. *Highlights in Science, Engineering and Technology*, 39:1164–1169.
- Charlie Winter and RV Gundur. 2022. Challenges in gaining ethical approval for sensitive digital social science studies. *International Journal of Social Research Methodology*, pages 1–16.
- Zimri S Yaseen, Mariah Hawes, Shira Barzilay, and Igor Galynker. 2019. Predictive validity of proposed diagnostic criteria for the suicide crisis syndrome: An acute presuicidal state. *Suicide and Life-Threatening Behavior*, 49(4):1124–1135.
- Michael Zimmer. 2020. “But the data is already public”: On the ethics of research in Facebook. In *The Ethics of Information Technologies*, pages 229–241. Routledge.
- Ayah Zirikly, Philip Resnik, Ozlem Uzuner, and Kristy Hollingshead. 2019. CLPsych 2019 shared task: Predicting the degree of suicide risk in Reddit posts. In *Proceedings of the Sixth Workshop on Computational Linguistics and Clinical Psychology*, pages 24–33.

## A Appendix Section

The following appendix section contains Fig 1. Topic word scores; Table 2. Topic frequency count; Fig 2. Hierarchical clustering; and Table 3. Illustrative composite synthetic quotes.

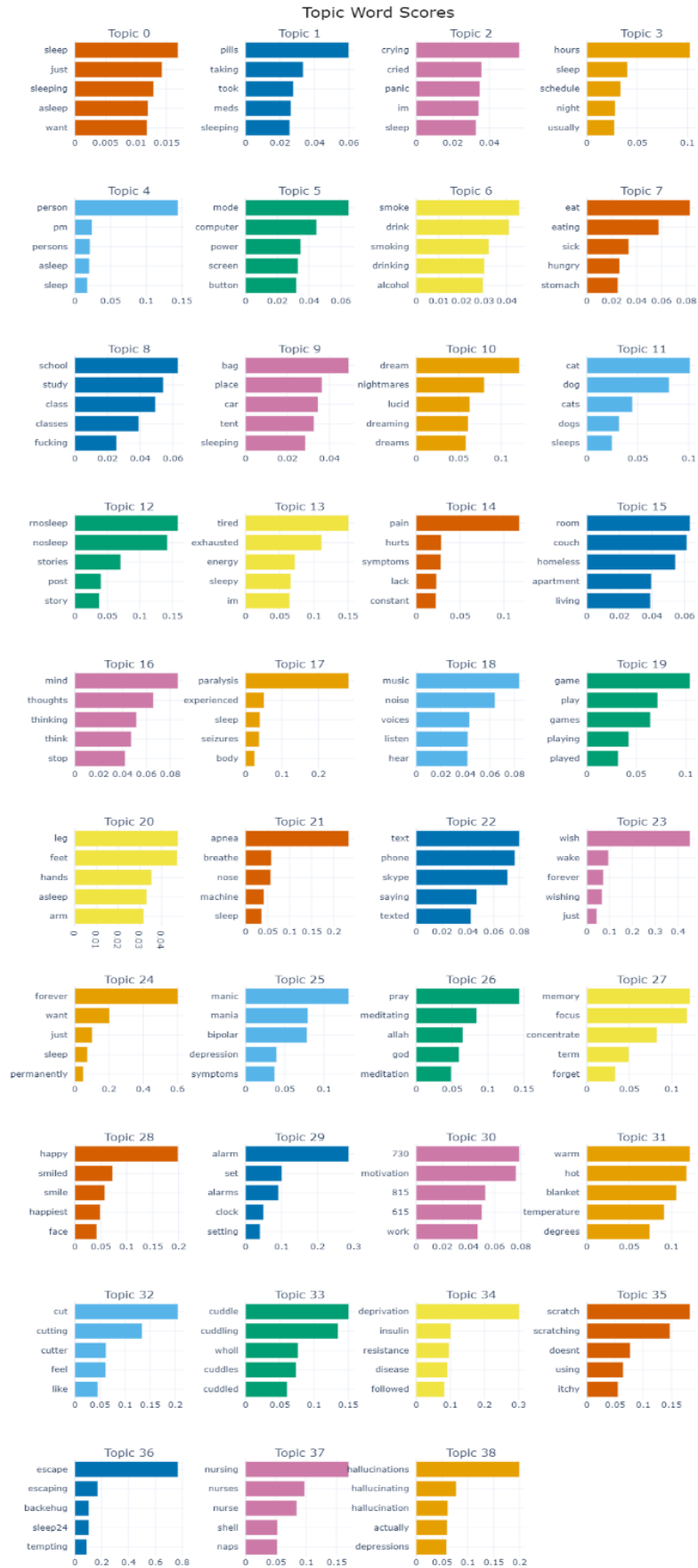


Figure 1: c-TF-IDF of key topic terms. TF-IDF (term frequency-inverse document frequency) is a statistical measure of the importance and relevancy of a word both to a specific document and across a corpus of documents).

<b>Topic</b>	<b>Frequency Count</b>	<b>Most important and relevant topic representative terms</b>
-1	9413	sleep_im_asleep_just
0	6130	sleep_just_sleeping_asleep
1	776	pills_taking_took_meds
2	620	crying_cried_panic_im
3	522	hours_sleep_schedule_night
4	383	person_pm_persons_asleep
5	366	mode_computer_power_screen
6	338	smoke_drink_smoking_drinking
7	320	eat_eating_sick_hungry
8	287	school_study_class_classes
9	270	bag_place_car_tent
10	263	dream_nightmares_lucid_dreaming
11	202	cat_dog_cats_dogs
12	178	nosleep_nosleep_stories_post
13	170	tired_exhausted_energy_sleepy
14	165	pain_hurts_symptoms_lack
15	159	room_couch_homeless_apartment
16	154	mind_thoughts_thinking_think
17	154	paralysis_experienced_sleep_seizures
18	150	music_noise_voices_listen
19	105	game_play_games_playing
20	101	leg_feet_hands_asleep
21	78	apnea_breathe_nose_machine
22	67	text_phone_skype_saying
23	60	wish_wake_forever_wishing
24	56	forever_want_just_sleep
25	48	manic_mania_bipolar_depression
26	44	pray_meditating_allah_god
27	43	memory_focus_concentrate_term
28	43	happy_smiled_smile_happiest
29	38	alarm_set_alarms_clock
30	26	730_motivation_815_615
31	25	warm_hot_blanket_temperature
32	22	cut_cutting_cutter_feel
33	22	cuddle_cuddling_wholl_cuddles
34	19	deprivation_insulin_resistance_disease
35	17	scratch_scratching_doesnt_using
36	16	escape_escaping_backehug_sleep24
37	15	nursing_nurses_nurse_shell
38	11	hallucinations_hallucinating_hallucination

Table 2: Frequency count of topics



## Hierarchical Clustering

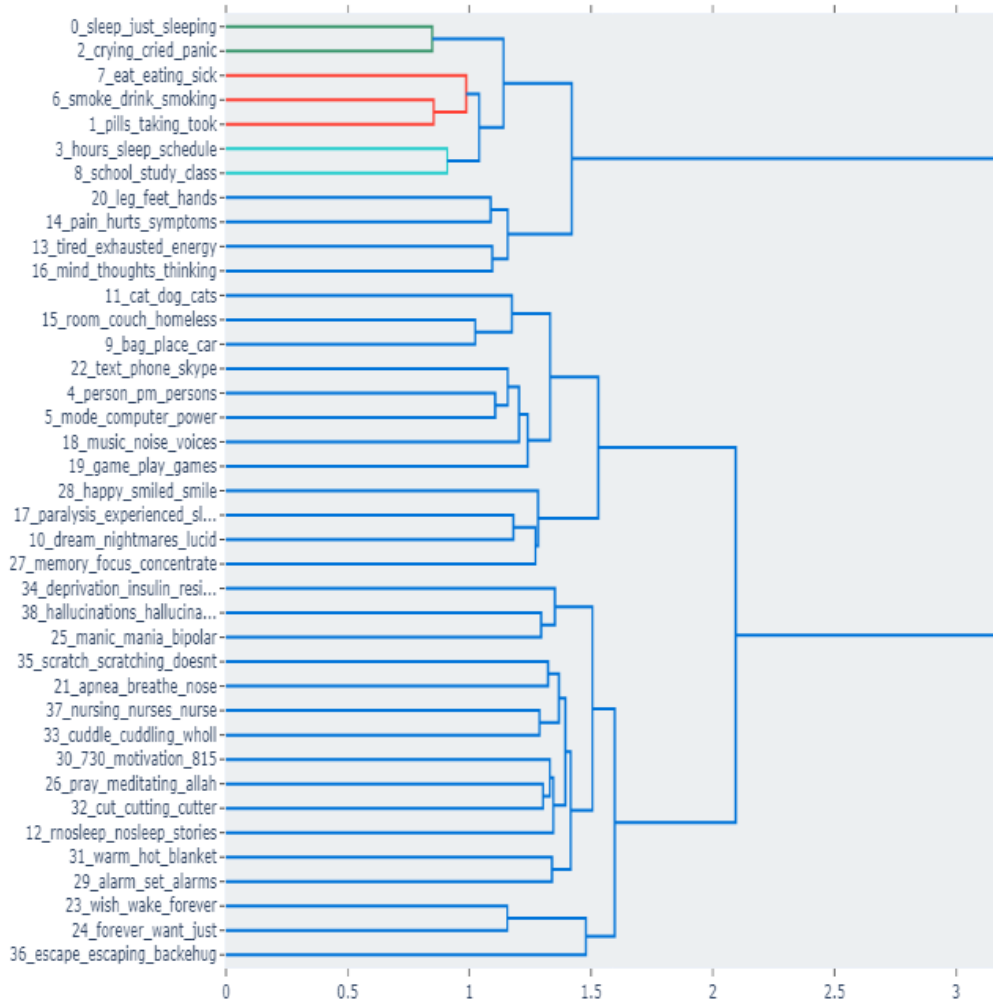


Figure 2: Hierarchical clustering of topics (utilizing `scipy.cluster.hierarchy` to create clusters and visualize how relate)

Exhausted, tired, emotionally and physically, just want to sleep forever, but never enough sleep
Every night I pray I die in my sleep, or someone ends the nightmare by putting a bullet in my head
Dreaming is a better form of life, particularly when can control with lucid dreaming
Calm the rampant thought storm, the panic, the mind at war that prevents me from sleeping
Tired of the loneliness, screwing up, and failure, dreams are my only escape and pleasure
Sleep use to be an escape but now a fear because of the thoughts, the sleep paralysis, and nightmares
I hide under the bed to feel safe, feeling entrapped and overwhelmed, and sleep the only thing to look forward to
Terrified of what the recesses of my mind and the darkness will conjure up
Just want to live in dreams and never wake up, escape the pain, grief, loneliness and shame
The worst thing I fear is the thoughts, the insomnia, and being alone at night
Wake up terrified from the nightmare at 4 am and too frightened to go back to sleep
The pain, the fear, the thoughts start from the moment I regain consciousness
I need to escape the suffocating thoughts and loneliness whether through drugs, sex, or sleep
Just want to knock myself out with pills or drugs so can stop feeling like shit and exhausted
I want someone to hold me and my trauma, and sadness until I fall asleep
Work, study, and bills I'm exhausted, sleep deprived, but I can't give up internet or risk sleeping in.
I lie in bed crying, I am a burden and a shame to my family, best thing for them is if I sleep forever

Table 3: Key illustrative composite synthetic quotes. A key issue in social media research is that verbatim quotes can not infrequently be traced back to the source, and with cross-referencing of other aspects of content or style, the online poster is potentially identified. These composite synthetic quotes, aim to illustrate the type of content that was the basis of the reflexive thematic analysis generated themes, without infringing the ethical undertaking not to use verbatim quotes.