

RESEARCH ARTICLE

# Identifying emerging mental illness utilizing search engine activity: A feasibility study

Michael L. Birnbaum<sup>1,2,3†\*</sup>, Hongyi Wen<sup>4†</sup>, Anna Van Meter<sup>1,2,3</sup>, Sindhu K. Ernala<sup>5</sup>, Asra F. Rizvi<sup>1,2,3</sup>, Elizabeth Arenare<sup>1,2,3</sup>, Deborah Estrin<sup>4</sup>, Munmun De Choudhury<sup>5</sup>, John M. Kane<sup>1,2,3</sup>

**1** The Zucker Hillside Hospital, Northwell Health, Glen Oaks, NY, United States of America, **2** The Feinstein Institute for Medical Research, Manhasset, NY, United States of America, **3** The Donald and Barbara Zucker School of Medicine at Hofstra/Northwell, Hempstead, NY, United States of America, **4** Cornell Tech, Cornell University, New York, NY, United States of America, **5** Georgia Institute of Technology, Atlanta, GA, United States of America

† These authors share first authorship on this work.

\* [Mbirnbaum@northwell.edu](mailto:Mbirnbaum@northwell.edu)



## OPEN ACCESS

**Citation:** Birnbaum ML, Wen H, Van Meter A, Ernala SK, Rizvi AF, Arenare E, et al. (2020) Identifying emerging mental illness utilizing search engine activity: A feasibility study. PLoS ONE 15(10): e0240820. <https://doi.org/10.1371/journal.pone.0240820>

**Editor:** Sinan Guloksuz, Department of Psychiatry and Neuropsychology, Maastricht University Medical Center, NETHERLANDS

**Received:** January 7, 2020

**Accepted:** October 4, 2020

**Published:** October 16, 2020

**Copyright:** © 2020 Birnbaum et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The datasets analyzed during the current study are not publicly available due to participant privacy and security concerns, including HIPAA regulations. The search archives and health records are not redistributable to researchers other than those engaged in the Institutional Review Board approved research collaborations. Data requests may be sent to the Northwell IRB at 1 (516) 465-1910 or [irb@northwell.edu](mailto:irb@northwell.edu).

## Abstract

Mental illness often emerges during the formative years of adolescence and young adult development and interferes with the establishment of healthy educational, vocational, and social foundations. Despite the severity of symptoms and decline in functioning, the time between illness onset and receiving appropriate care can be lengthy. A method by which to objectively identify early signs of emerging psychiatric symptoms could improve early intervention strategies. We analyzed a total of 405,523 search queries from 105 individuals with schizophrenia spectrum disorders (SSD, N = 36), non-psychotic mood disorders (MD, N = 38) and healthy volunteers (HV, N = 31) utilizing one year's worth of data prior to the first psychiatric hospitalization. Across 52 weeks, we found significant differences in the timing ( $p < 0.05$ ) and frequency ( $p < 0.001$ ) of searches between individuals with SSD and MD compared to HV up to a year in advance of the first psychiatric hospitalization. We additionally identified significant linguistic differences in search content among the three groups including use of words related to sadness and perception, use of first and second person pronouns, and use of punctuation (all  $p < 0.05$ ). In the weeks before hospitalization, both participants with SSD and MD displayed significant shifts in search timing ( $p < 0.05$ ), and participants with SSD displayed significant shifts in search content ( $p < 0.05$ ). Our findings demonstrate promise for utilizing personal patterns of online search activity to inform clinical care.

## Introduction

The consequences of untreated psychiatric illness can be devastating [1–3]. Behavioral health disorders often present during the formative years of adolescent and young adult development and interfere with the establishment of social, educational, and vocational foundations [4]. While early intervention services have demonstrated the potential to improve outcomes, symptoms often remain unrecognized and untreated for years before receiving effective care

**Funding:** The authors received no specific funding for this work.

**Competing interests:** The authors have declared that no competing interests exist.

[5–8]. Novel screening strategies, supported by technological innovation, are critical to achieving the goal of early identification and treatment.

The emergence of serious mental illnesses, such as schizophrenia and bipolar disorder, are often preceded by periods of anxiety, mood lability, sleep pattern irregularity, trouble concentrating, social isolation, strained interactions with others, and attenuated/subthreshold psychotic and manic experiences [9, 10]. Despite the decline in functioning and established deleterious impact of untreated symptoms, an effective method by which to screen and educate vulnerable individuals has not been established [11, 12]. Clinical interview, assessment scales, patient self-report, and family observation remain the primary sources for assessing early warning signs and are limited by reliance on direct and timely contact with trained professionals, as well as accurate and insightful patient and family recall. These standard approaches to clinical assessment do not allow for objective monitoring of psychiatric symptom emergence and typically do not occur with enough frequency, or at the necessary level of detail, to detect subtle, sub-clinical, and burgeoning symptoms. Early, precise, and noninvasive identification of psychiatric symptom emergence could facilitate the initiation of personalized and proactive intervention strategies.

At the same time, Google Search has emerged as one of the world's most popular websites, supporting over 660 million daily visitors, and managing over three billion searches daily [13]. Searching online has become the primary resource for youth seeking out mental health related information [14]. This is especially true for stigmatized illnesses such as schizophrenia, as the Internet provides a comfortable and anonymous setting to gather information about symptoms and treatment options [15]. Previous reports have demonstrated that adolescents and young adults with emerging symptoms of psychiatric disorders utilize the Internet first, and most frequently, to gather information prior to receiving psychiatric care, and that they are more likely to search online for information than to discuss their experiences with peers, family, physicians, and mental health clinicians [16–18]. Performing an Internet search may therefore represent one of the first proactive steps towards treatment initiation and could provide a valuable opportunity to impact help-seeking behavior.

Prior work in machine learning has highlighted opportunities to utilize large scale anonymized online search activity to detect content and patterns associated with the emergence and progression of medical illnesses including lung cancer, pancreatic cancer, and Parkinson's disease [19–21]. These initiatives aim to inform the development of a new generation of digital tools designed to assist in the screening and early identification of individuals developing medical health conditions. Similar computational methods have identified associations between social media activity and behavioral health [22–29]. Few studies to date, however, have explored the link between search activity and psychiatric illness, beyond retrospective self-report [30]. Furthermore, while promising, internet activity research to date has been limited by the fact that it has been conducted nearly exclusively using search data from anonymous individuals who self-report a diagnosis online, and has yet to be carried out in real world clinical settings, using participant-contributed search data, with clinically validated symptoms and diagnoses.

This study aimed to explore the feasibility of utilizing online search archives as a tool to identify emerging psychiatric symptoms. This knowledge would support the development of resources designed to inform screening procedures for individuals with emerging mental illness earlier along their trajectory to care. We hypothesized that significant differences in the *content, timing, and frequency* of online activity would differentiate participants with schizophrenia spectrum disorders (SSD) from those with mood disorders (MD) and healthy volunteers (HV). Additionally, we hypothesized that significant changes in the content and behavioral patterns of search activity would exist within individuals with SSD and MD in the

period of time closest to their first hospitalization consistent with escalating psychiatric symptoms during that time period.

## Materials and methods

Participants between the ages of 15 and 35 years, diagnosed with a schizophrenia spectrum disorder or a non-psychotic mood disorder, were recruited from The Zucker Hillside Hospital/Northwell Health inpatient and outpatient psychiatric departments. Most participants with SSD were recruited from the Early Treatment Program (ETP), Zucker Hillside's specialized early psychosis intervention clinic. Additional participants (3) were recruited from a collaborating psychiatric clinic located in East Lansing, Michigan. Healthy volunteers who had already been screened for prior studies were recruited. Additional HV's ( $n = 10$ ) were recruited from the University of North Carolina. Recruitment occurred between March 2016 and December 2018. Written informed consent was obtained for adult participants and legal guardians of participants under 18 years of age. Assent was obtained for participating minors. Participants were fully informed of the potential risks, benefits, and alternative options available, as well as strategies to mitigate risks. Decisional capacity to consent was determined through clinical assessment, as well as via completion of a short quiz, designed to assess one's understanding of research procedures, conducted prior to consenting to participate. The study was approved by the Institutional Review Board (IRB) of Northwell Health (the coordinating institution) as well as local IRBs at participating sites.

Participation involved a single study visit. Participants were asked to export their search archive by logging on to their Google account to request their search history. Archives include all historical search activity including the content and timing of search queries. Diagnoses and dates for the first psychiatric hospitalization were obtained through participants' medical records.

Given the goal of identifying changes in search activity associated with escalating psychiatric symptoms, 52 weeks' worth of search data prior to the first psychiatric hospitalization was extracted from each participant, operating under the expectation that at some point during that year, psychiatric signs and symptoms emerged and progressed to the point of necessitating inpatient intervention. One year was selected as it represents a period of time long enough to establish a baseline level of search activity, and to identify changes in the weeks closest to hospitalization. For HV (who were never hospitalized), the midpoint of the first hospitalization dates across all patients ( $N = 74$ ) was utilized to mitigate the potential temporal effects on search patterns, such as functional changes in the search platform and search data logging systems over time. This resulted in using November 9, 2015 as the anchor date for healthy participants in our dataset.

Our analysis consisted of (1) between-group comparisons among SSD, MD, and HV to examine group-level differences, and (2) within group comparisons by comparing a period of "relative health" (6 month furthest away from hospitalization) to periods of "relative illness", closest to the date of the first psychiatric hospitalization.

Both sets of comparisons were conducted on the frequency, timing, and content of searches. We extracted search frequency and timing distributions from the meta data (i.e. timestamps). For search content, we used Linguistic Inquiry and Word Count (LIWC), a well validated language analytic tool, which extracts 93 variables pertaining to word usage, known to be associated with emotion, mood, and behavior [31, 32]. Given the number of comparisons tested, we implemented the two-stage Benjamini and Hochberg [33] procedure to control the false discovery rate (FDR). Specifically, we used the implementation from the statsmodels Python library [34] and set the family-wise error rate to be 0.05.

## Data preprocessing

A total of 132 (44 SSD, 41 MD, 47 HV) search archives were available for analysis. Participants with 30 or more weeks of zero search activity during the 52-week period were excluded ( $n = 27$ ). The final dataset consisted of 405,523 searches across 105 participants. Demographic information of included participants is shown in Table 1.

## Results

### Between-group differences in search (frequency, timing, and content)

**Frequency of search activity.** Across 52-weeks (Fig 1), HV showed significantly higher search frequency on average compared to both SSD (Post-hoc Tukey:  $T = 19.51$ ,  $p = 0.001$ ) and MD (Post-hoc Tukey:  $T = 16.76$ ,  $p = 0.001$ ). HV also showed significantly higher variability of search frequency compared to MD ( $T = 2.83$ ,  $p = 0.006$ ) across the 52-weeks (averaged standard deviations across weeks: HV = 50.74, MD = 30.12, SSD = 36.36). Education, sex, and age were not associated with search frequency among MD and HV participants. Among SSD participants, those who completed high school ( $n = 25$ ) searched more often than those who did not ( $n = 11$ ), and relatedly, young adults with SSD, 20 years and older ( $n = 28$ ), searched more than adolescents ( $n = 8$ ).

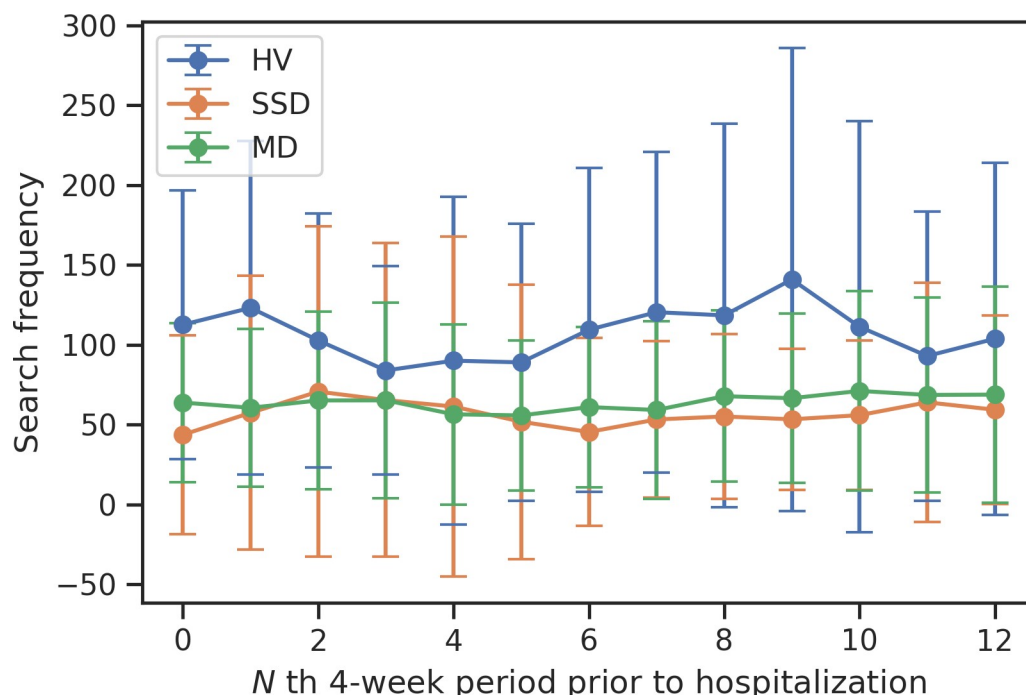
**Timing of search activity.** Over 52 weeks (Fig 2), we found that SSD participants search significantly more during the 12am-6am period ( $T = 2.24$ ,  $p = 0.029$ ) compared to HV. Additionally, MD participants searched significantly less than HV from 12pm-6pm ( $T = -2.20$ ,  $p = 0.03$ ) and significantly more than SSD from 6pm-12am ( $T = 2.48$ ,  $p = 0.015$ ). Education, sex, and age were not associated with search timing among MD, SSD, and HV participants.

**Content of search activity.** Over 52-weeks (Table 2), we identified several linguistic differences in search content across groups. Participants with SSD were significantly more likely

Table 1. Participant demographics.

	SSD	MD	HV	Full Sample
N	36	38	31	105
	Mean (SD)			
Age	23.11 (3.3)	19.48 (3.1)	25.72 (4.8)	23.12 (4.2)
Years of Education	13.58 (1.8)	13.34 (2.1)	16.41 (1.9)	14.29(2.3)
	n (%)			
Male	22 (61)	10 (26)	11 (35)	43 (41)
Race/Ethnicity				
African American/Black	16 (44)	7 (18.4)	5 (16.1)	28 (27)
Asian	5 (13.9)	5 (13.2)	6 (19.4)	16 (15)
Caucasian	12 (33.3)	11 (28.9)	18 (58.1)	47 (45)
Mixed race/Other	3 (8.3)	9 (23.7)	2 (6.4)	14 (13)
Hispanic	9 (25)	11 (28.9)	0 (0)	20 (19)
Diagnosis				
Schizophrenia	16 (15)	0 (0)	0 (0)	16 (15)
Schizophreniform	8 (8)	0 (0)	0 (0)	8 (8)
Schizoaffective	1 (1)	0 (0)	0 (0)	1 (0)
Brief Psychotic Disorder	2 (2)	0 (0)	0 (0)	2 (2)
Unspecified SSD	9 (9)	0 (0)	0 (0)	9 (9)
Bipolar I Disorder	0 (0)	5 (5)	0 (0)	5 (5)
Major Depressive Disorder	0 (0)	33 (32)	0 (0)	33 (31)

<https://doi.org/10.1371/journal.pone.0240820.t001>



**Fig 1.** Search frequency across groups over 52 weeks.

<https://doi.org/10.1371/journal.pone.0240820.g001>

to search using words related to perception ( $T = 3.08$ ,  $p = 0.025$ ) and use first ( $T = 3.01$ ,  $p = 0.03$ ) and second person pronouns ( $T = 3.45$ ,  $p = 0.011$ ) compared to HV. Participants with MD were significantly more likely to search using words related to negative emotions ( $T = 2.94$ ,  $p = 0.028$ ), sadness ( $T = 3.01$ ,  $p = 0.026$ ), and death ( $T = 2.71$ ,  $p = 0.046$ ), and use first ( $T = 3.41$ ,  $p = 0.010$ ) and second ( $T = 3.22$ ,  $p = 0.016$ ) person pronouns compared to HV. HV were significantly more likely to search using more words compared to SSD ( $T = 3.57$ ,  $p = 0.009$ ) and MD ( $T = 3.18$ ,  $p = 0.018$ ), use more punctuation compared to SSD ( $T = 3.53$ ,  $p = 0.009$ ) and MD ( $T = 3.27$ ,  $p = 0.015$ ), and use common online abbreviations (i.e., b/c for “because”) compared to SSD ( $T = 3.42$ ,  $p = 0.011$ ).

### Within-group differences in search (frequency, timing, and content)

To explore within group differences in search frequency, timing, and content, search data was aggregated and averaged over one-week intervals. The periods of time within 6 months (24 weeks) closest to hospitalization were defined as periods of “relative illness” as we would expect symptoms to be most prominent during this time, culminating in hospitalization. These periods were compared to periods of “relative health,” which consisted of data from the six months (25–52 weeks) furthest away from hospitalization.

**Frequency of search activity.** No significant differences in search frequency were found between periods of relative illness and periods of relative health in all three groups using repeated measures ANOVA and paired t-tests.

**Timing of search activity.** Significant shifts were identified in the timing of search activity in participants with MD and SSD closer to hospitalization (Figs 3 and 4). Compared to periods of relative health, participants with MD searched significantly less ( $T = -3.19$ ,  $p = 0.003$ ) during the morning hours (6am–12pm) during periods of relative illness. Compared to periods of relative health, participants with SSD search significantly less ( $T = -2.30$ ,  $p = 0.03$ ) during the early morning hours (12am–6am) during periods of relative illness.

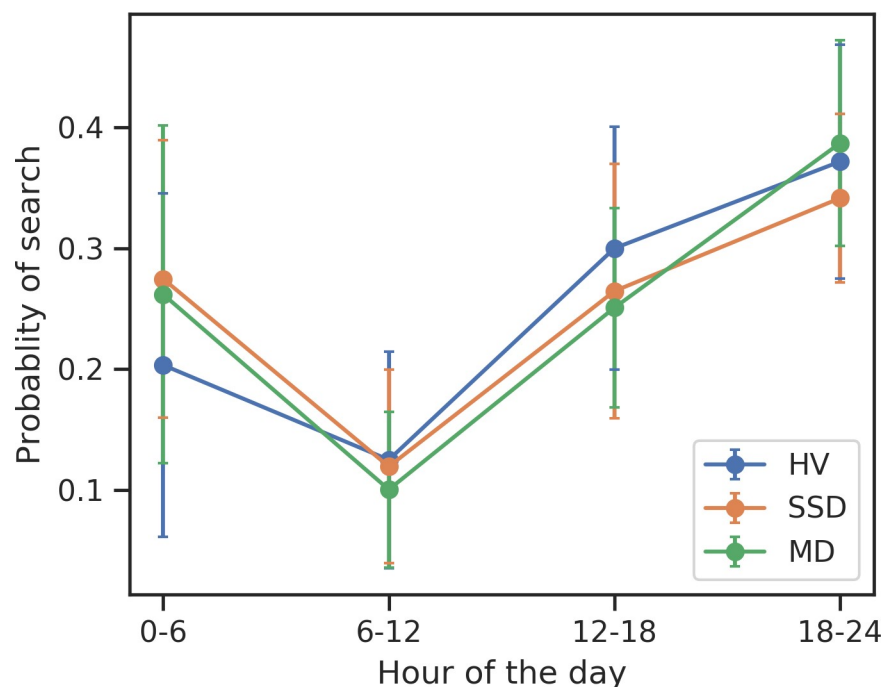


Fig 2. Search timing across groups over 52 weeks.

<https://doi.org/10.1371/journal.pone.0240820.g002>

**Content of search activity.** We identified several significant linguistic shifts in search content among participants with SSD prior to the first psychiatric hospitalization (Table 3). Participants with SSD were less likely to use punctuation ( $T = -4.13, p = 0.0025$ ), less likely to search for terms related to “seeing” ( $T = -3.79, p = 0.01$ ), “anger” ( $T = -3.47, p = 0.023$ ), “negative emotions” ( $T = -3.15, p = 0.04$ ), “perception” ( $T = -3.30, p = 0.03$ ), and “death” ( $T = -3.54, p = 0.04$ ), in the 12 weeks prior to hospitalization compared to periods of relative health (Figs 5–7). No significant shifts in search content were identified among the MD participants.

Table 2. Linguistic differences in search content across 52 weeks.

HV > MD	HV > SSD	MD > HV	MD > SSD	SSD > MD	SSD > HV
Word count ( $p = 0.018$ )	Word count ( $p = 0.009$ )	Authentic ( $p = 0.010$ )	Filler	Semi-Colon	Authentic ( $p = 0.009$ )
Analytic ( $p = 0.006$ )	Analytic ( $p = 0.009$ )	Function ( $p = 0.006$ )	( $p = 0.045$ )	( $p = 0.045$ )	Total pronoun
Words per sentence	Words per sentence	Total pronoun ( $p = 0.006$ )			( $p = 0.009$ )
( $p = 0.028$ )	( $p = 0.013$ )	Personal pronouns			You ( $p = 0.011$ )
Punctuation ( $p = 0.015$ )	We ( $p = 0.023$ )	( $p = 0.006$ )			I ( $p = 0.030$ )
Period ( $p = 0.018$ )	Informal ( $p = 0.013$ )	I ( $p = 0.010$ )			Perception ( $p = 0.025$ )
Dash ( $p = 0.023$ )	Net speak ( $p = 0.011$ )	You ( $p = 0.016$ )			Focus Present ( $p = 0.009$ )
	Punctuation ( $p = 0.009$ )	Prep ( $p = 0.040$ )			
	Period ( $p = 0.009$ )	Aux verb ( $p = 0.006$ )			
	Dash ( $p = 0.021$ )	Adverb ( $p = 0.006$ )			
		Negate ( $p = 0.006$ )			
		Verb ( $p = 0.006$ )			
		Interrogation ( $p = 0.027$ )			
		Quant ( $p = 0.045$ )			
		Neg emotion ( $p = 0.028$ )			
		Sad ( $p = 0.026$ )			
		Cog process ( $p = 0.027$ )			
		Cause ( $p = 0.046$ )			
		Focus present ( $p = 0.006$ )			
		Death ( $p = 0.046$ )			

<https://doi.org/10.1371/journal.pone.0240820.t002>

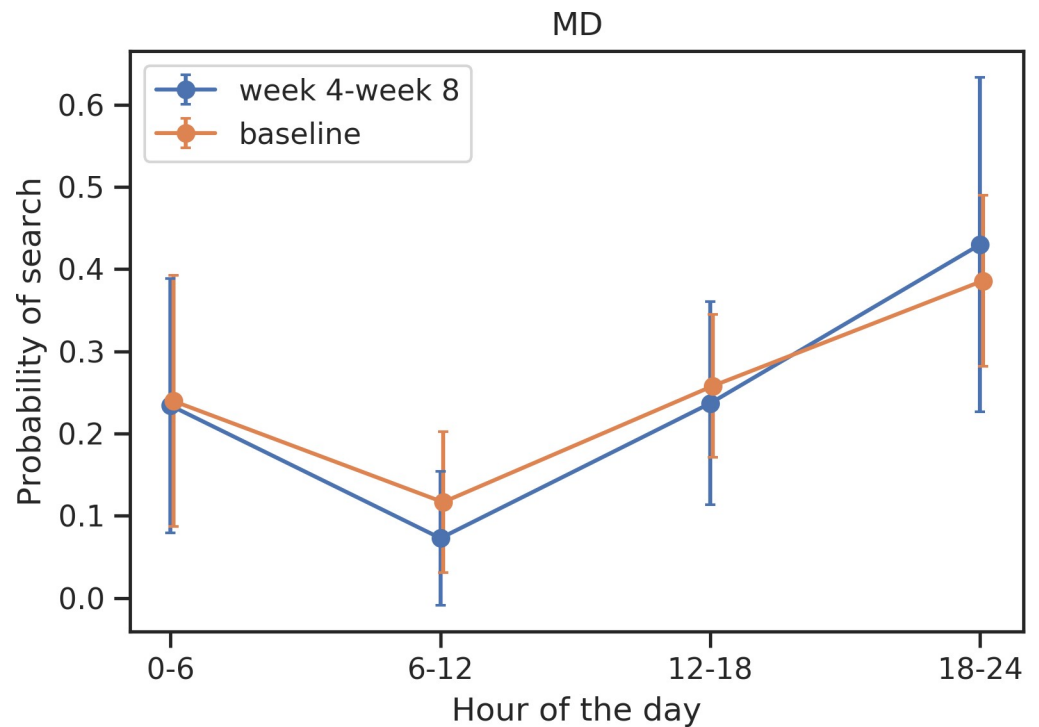


Fig 3. Shifts in timing of search activity across 24 hours in participants with MD.

<https://doi.org/10.1371/journal.pone.0240820.g003>

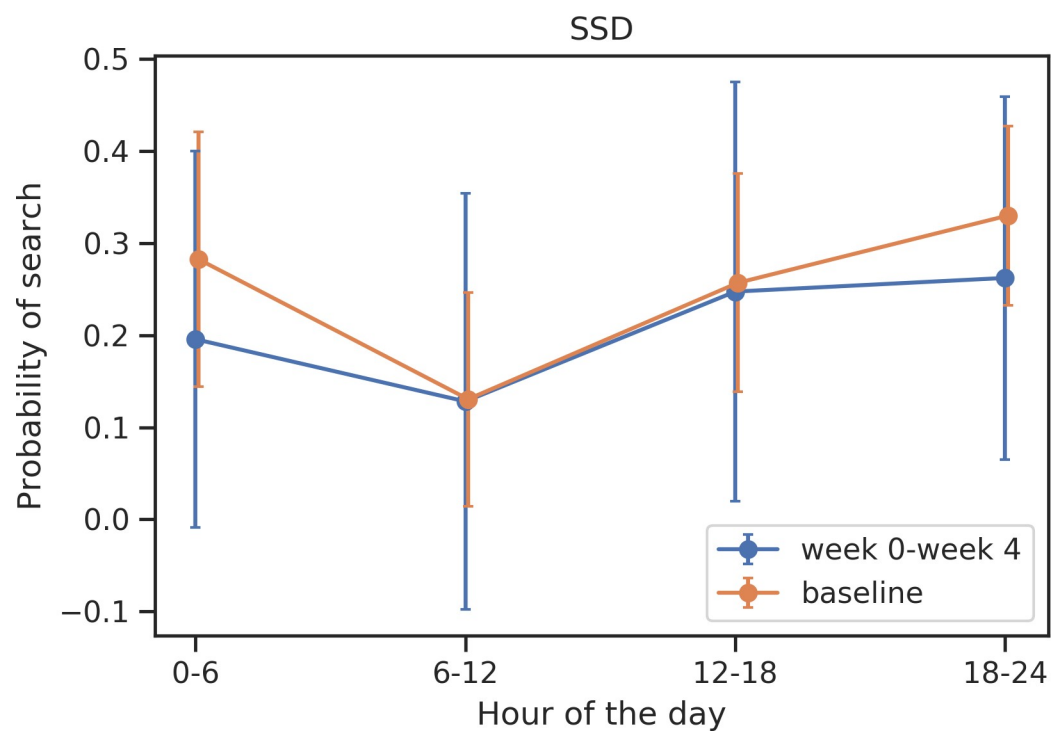


Fig 4. Shifts in timing of search activity across 24 hours in participants with SSD.

<https://doi.org/10.1371/journal.pone.0240820.g004>



Table 3. Within group changes in content for SSD (comparing periods of relative illness to periods of relative health).

	Weeks prior to hospitalization		
	Week 0–4	Week 4–8	Week 8–12
Content of searches (Relative health > Relative illness)	Quote ( $p = 0.025$ ) Comma ( $p = 0.000008$ ) Punctuation ( $p = 0.0025$ )	Neg emotion ( $p = 0.04$ ) Anger ( $p = 0.025$ ) Perception ( $p = 0.03$ ) See ( $p = 0.01$ )	Death ( $p = 0.04$ )

<https://doi.org/10.1371/journal.pone.0240820.t003>

## Discussion

In this study, we explored the potential for online search activity to serve as a tool to identify emerging behavioral health disorders. Our results suggest significant differences exist in the timing, frequency, and content of search activity a year in advance of the first psychiatric hospitalization for participants with SSD and MD, when compared to HV. Furthermore, in the weeks closest to the date of the first hospitalization, significant shifts in language occurred in participants with SSD, and significant shifts in timing occurred in individuals with SSD and MD. While Google data alone is not meant to diagnose psychiatric conditions, results demonstrate the potential for online search activity to be used in conjunction with clinical information to inform clinical decision making. Similar to the way a physician might use an x-ray or blood test to inform health status, search data may one day serve as a viable screening tool to better gauge risk factors associated with the later development of psychiatric conditions. Identifying emerging psychiatric symptoms early, before they have an opportunity to escalate to the point of necessitating hospitalization, is our best chance at transforming trajectories to care and improving behavioral healthcare experiences and outcomes for patients.

Individuals with SSD and MD demonstrated significantly fewer searches in the year leading up to the first psychiatric hospitalization as compared to HV, who searched over twice as much. Counter to our hypothesis, no significant changes in search frequency were noted in participants with SSD and MD as psychiatric symptoms escalated necessitating a psychiatric hospitalization. Decreased search activity may be related to very early budding psychiatric symptoms including reduced motivation, increased fatigue, or decreased interest and engagement with one's environment [35–37]. Additionally, individuals with SSD are known to experience cognitive deficits early in the course of illness development, which may contribute to reduced search activity [38, 39]. These subtle changes may occur well in advance of the first psychiatric hospitalization. In contrast, reduced search activity may represent a longstanding risk factor contributing to the later emergence of a psychiatric disorder. To address these questions, future research will need to extract search data several years in advance of the first psychiatric hospitalization, and to prospectively collect symptom rating scales in individuals earlier along the course of illness development.

Compared to HV, participants with SSD and MD searched at different times throughout the day. Temporal differences date back at least a year in advance of the first psychiatric hospitalization. Sleep disruption is a common experience for people with psychiatric disorders [40, 41], and many individuals with MD show circadian shifting [42], which results in a preference for being awake/active late at night. Precisely when sleep disturbances begin, however, is less well understood. According to our data, sleep dysfunction appears to already be present well in advance of the first psychiatric hospitalization and significant alterations in search timing occurred in both patient populations closer to the date of the first psychiatric hospitalization. As with search frequency, it remains unclear if different temporal patterns represent a change from baseline activity due to emerging psychiatric symptoms, or rather a persistent irregularity in sleep contributing to the later development of a psychiatric disorder. In either circumstance,



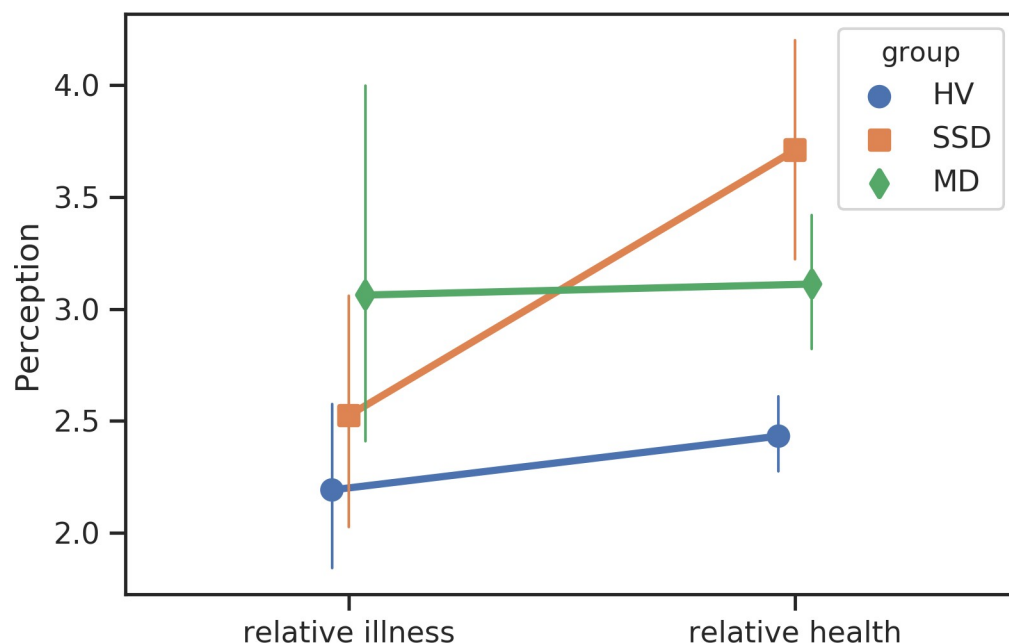


Fig 5. Changes in search content corresponding to “perception”.

<https://doi.org/10.1371/journal.pone.0240820.g005>

extracting online search activity in youth presenting with sleep disturbance may one day serve as useful collateral information to predict the risk of psychiatric illness development.

Linguistic analysis of search terms identified significant differences in search content over 52-weeks before the first psychiatric hospitalization. Compared to HV, search terms among participants with SSD and MD demonstrated a greater emphasis on sadness, and perception,

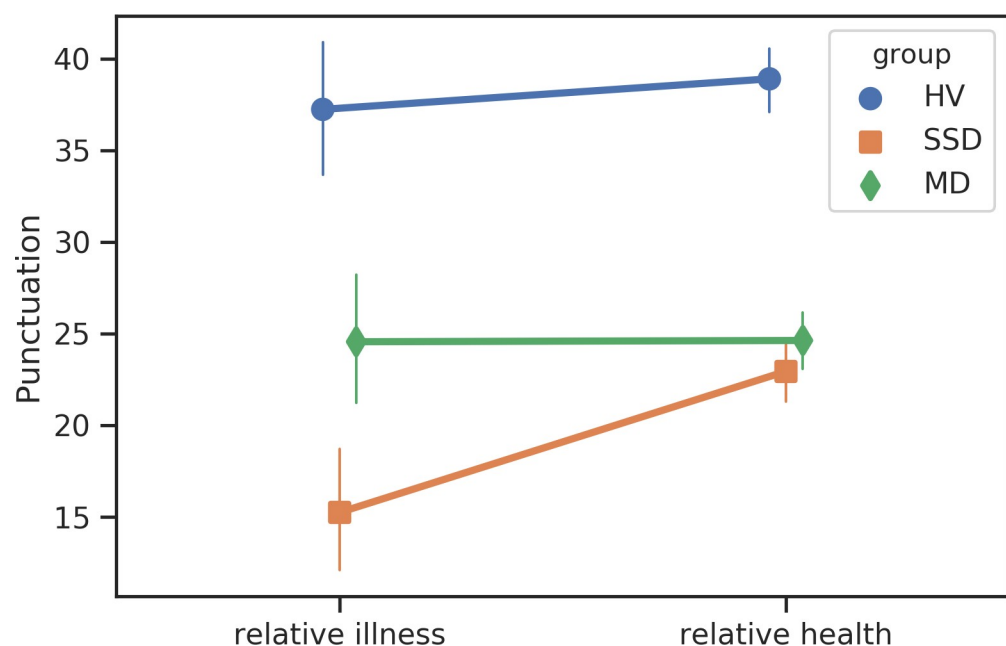


Fig 6. Changes in search content corresponding to “punctuation”.

<https://doi.org/10.1371/journal.pone.0240820.g006>

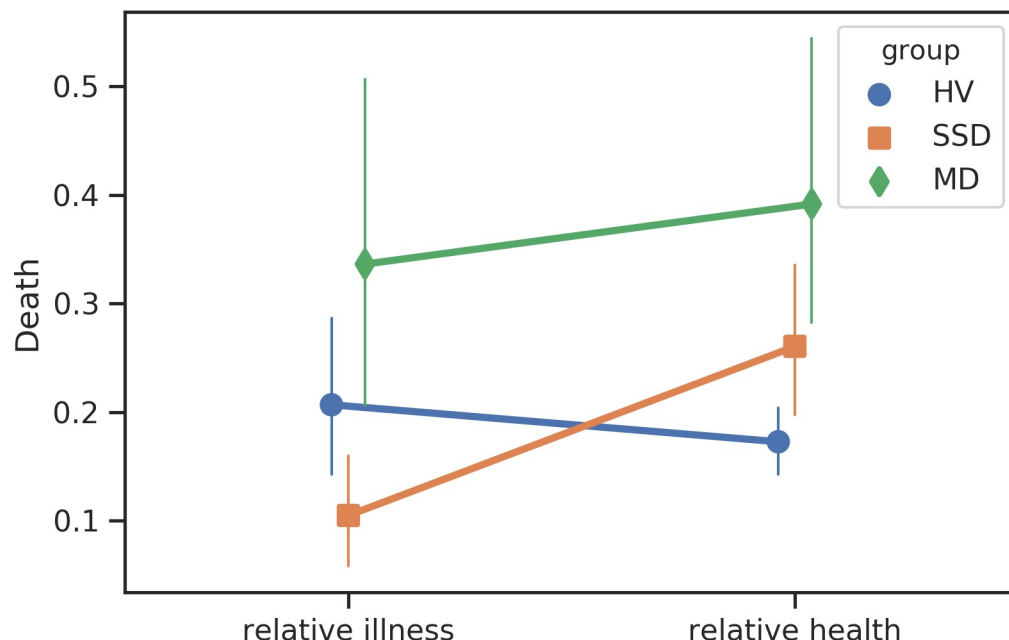


Fig 7. Changes in search content corresponding to “death”.

<https://doi.org/10.1371/journal.pone.0240820.g007>

as well as first and second person pronouns. Several linguistic differences existed well in advance of the first psychiatric hospitalization and it is possible that certain linguistic features represent a state rather than a trait marker of mental illness. Analyzing word choice could therefore help to identify people at higher risk of SSD or MD prior to the emergence of clinically significant symptoms. As symptoms progressed, closer to the date of hospitalization, the content of searches changed significantly among individuals with SSD, but not MD. These linguistic changes may reflect shifting interests, changing mood, preoccupations, social functioning, and other domains known to accompany psychotic illness emergence [4]. In contrast to prior work exploring changes in language use on social media associated with relapse [29], participants with SSD were less likely to search for content related to perceptions, anger, and negative emotions. This may be related to differences in how people compose searches, which are private and generally intended to find information, versus social media posts, which are public and may be more likely to be communicating information.

Prior research in linguistic analysis has identified significant differences at the word level in the use of certain word categories, as well as at the sentence level in terms of semantic density, coherence, and/or content, both in individuals at risk for developing psychotic disorders as well as those with established SSD and MD [43–58]. Language extracted from online search activity is distinct due to its short sentence structure and the unique nature of the online search platform. As we continue to identify linguistic associations with psychiatric illness, the source of the language data must be taken into consideration. Future work is needed to better understand the clinical correlates of changing online search content and to identify the point in illness progression at which linguistic shifts emerge, in order to make the best clinical use of this information. Additionally, further analysis is needed to identify how language varies depending on the platform (Facebook vs Google, for example) used, and which has greater clinical utility.

Several noteworthy limitations should be mentioned. First, our sample size is relatively small and limits the generalizability of these findings. Second, the majority of our participants

had medical record documentation that began with the first psychiatric hospitalization, making it challenging to know what symptoms were present and for how long prior to hospitalization. Given that many individuals report extended periods of untreated illness or comorbid psychiatric conditions prior to receiving clinical attention [5–8], it is possible that more data, beyond 52-weeks, is needed to identify shifts in search activity associated with illness progression. Additionally, future studies should consider monitoring participants prospectively and leveraging rating scales to more accurately explore how individual fluctuating interests and psychiatric symptoms impact search behaviors over time. Third, the fact that some participants searched more than others, may also impact results as there were large differences in the amount of extractable data across participants. While we do not anticipate that Google data alone will ever be sensitive or specific enough to establish a particular diagnosis, important questions for future research will be how much search data is necessary to make a reliable clinical prediction and how individual characteristics influence search behavior within a diagnostic group. Finally, eligibility criteria ranged from 15 to 35 years to reflect the inclusion criteria of the Early Treatment Program, however adolescents may engage with the Internet in a distinct manner compared to young adults and future initiatives will need to consider the impact of age as well as other demographic characteristics, such as sex, and education level on search activity.

Search patterns hold promise for gathering objective, non-invasive, and easily accessed, indicators of psychiatric symptom emergence. Utilizing online activity as collateral behavioral health information would represent a major advancement in efforts to capitalize on objective digital data to improve mental health screening. This would be a significant step forward for psychiatry, which has historically been limited by its reliance on self-reported data. However, how to effectively and ethically incorporate personalized patterns of online activity into public health initiatives and clinical workflow are critical questions [59]. The data utilized in the current study were obtained from consenting participants who were fully informed of the risks and benefits of participation. Furthermore, the data were extracted and analyzed locally at Northwell Health and remained entirely within a HIPAA compliant secure database to ensure the privacy of our participants. Nonetheless, this field of research evokes a host of challenging questions and concerns related to ethics, privacy, consent, and clinical responsibility. Interdisciplinary teams of researchers, clinicians, and patients must continue to work together on identifying and solving these important ethical dilemmas. Importantly, investigators must develop standards to protect the confidentiality and the rights of this sensitive population to avoid misuse of personal information and ensure that the data and the technologies are used in the service of positive outcomes for clinicians and the patients they treat.

## Author Contributions

**Conceptualization:** Michael L. Birnbaum, Anna Van Meter, Sindhu K. Ernala, Asra F. Rizvi, Elizabeth Arenare, Munmun De Choudhury, John M. Kane.

**Data curation:** Michael L. Birnbaum, Sindhu K. Ernala, Asra F. Rizvi, Elizabeth Arenare.

**Formal analysis:** Michael L. Birnbaum, Hongyi Wen, Sindhu K. Ernala, Munmun De Choudhury.

**Investigation:** Michael L. Birnbaum, Hongyi Wen.

**Methodology:** Michael L. Birnbaum, Hongyi Wen, Anna Van Meter, Asra F. Rizvi.

**Project administration:** Michael L. Birnbaum.

**Resources:** Michael L. Birnbaum.

**Software:** Hongyi Wen, Sindhu K. Ernala.

**Supervision:** Michael L. Birnbaum, Deborah Estrin, Munmun De Choudhury, John M. Kane.

**Writing – original draft:** Michael L. Birnbaum, Hongyi Wen, Anna Van Meter, Sindhu K. Ernala.

**Writing – review & editing:** Michael L. Birnbaum, Hongyi Wen, Anna Van Meter, Sindhu K. Ernala, Munmun De Choudhury, John M. Kane.

## References

1. Marshall M, Lewis S, Lockwood A, Drake R, Jones P, Croudace T. Association between duration of untreated psychosis and outcome in cohorts of first-episode patients: a systematic review. *Arch Gen Psychiatry*. 2005 Sep; 62(9):975–83. <https://doi.org/10.1001/archpsyc.62.9.975> PMID: 16143729
2. Perkins DO, Gu H, Boteva K, Lieberman JA. Relationship between duration of untreated psychosis and outcome in first-episode schizophrenia: a critical review and meta-analysis. *Am J Psychiatry*. 2005 Oct; 162(10):1785–804. Review. <https://doi.org/10.1176/appi.ajp.162.10.1785> PMID: 16199825
3. Ghio L, Gotelli S, Marcenaro M, Amore M, Natta W. Duration of untreated illness and outcomes in unipolar depression: a systematic review and meta-analysis. *J Affect Disord*. 2014; 152–154:45–51. <https://doi.org/10.1016/j.jad.2013.10.002> PMID: 24183486
4. American Psychiatric Association. Diagnostic and statistical manual of mental disorders. 5th ed. Washington DC: APA Publishing; 2013
5. Addington J, Heinssen RK, Robinson DG, Schooler NR, Marcy P, Brunette MF, et al. Duration of Untreated Psychosis in Community Treatment Settings in the United States. *Psychiatr Serv*. 2015 Jul; 66(7):753–6. <https://doi.org/10.1176/appi.ps.201400124> Epub 2015 Jan 15. PMID: 25588418
6. Kane JM, Robinson DG, Schooler NR, Mueser KT, Penn DL, Rosenheck RA, et al. Comprehensive Versus Usual Community Care for First-Episode Psychosis: 2-Year Outcomes From the NIMH RAISE Early Treatment Program. *Am J Psychiatry*. 2016 Apr 1; 173(4):362–72. <https://doi.org/10.1176/appi.ajp.2015.15050632> Epub 2015 Oct 20. PMID: 26481174; PubMed Central PMCID: PMC4981493
7. Vieta E, Salagre E, Grande I, Carvalho AF, Fernandes BS, Berk M, et al. Early intervention in bipolar disorder. *Am J Psychiatry*. 2018 Jan 24; 175(5):411–26. <https://doi.org/10.1176/appi.ajp.2017.17090972> PMID: 29361850
8. Duffy A, Jones S, Goodday S, Bentall R. Candidate Risks Indicators for Bipolar Disorder: Early Intervention Opportunities in High-Risk Youth. *Int J Neuropsychopharmacol*. 2015; 19(1):pyv071. Published 2015 Jun 25. <https://doi.org/10.1093/ijnp/pyv071> PMID: 26116493
9. Birchwood M, Spencer E, McGovern D. Schizophrenia: early warning signs. *Adv Psychiatr Treat*. 2000 Mar; 6(2):93–101
10. Van Meter AR, Burke C, Youngstrom EA, Faedda GL, Correll CU. The Bipolar Prodrome: Meta-Analysis of Symptom Prevalence Prior to Initial or Recurrent Mood Episodes. *J Am Acad Child Adolesc Psychiatry*. 2016; 55(7):543–555. <https://doi.org/10.1016/j.jaac.2016.04.017> PMID: 27343882
11. Lakkis NA, Mahmassani DM. Screening instruments for depression in primary care: a concise review for clinicians. *Postgrad Med*. 2015; 127(1):99–106. <https://doi.org/10.1080/00325481.2015.992721> PMID: 25526224
12. Addington J, Stowkowy J, Weiser M. Screening tools for clinical high risk for psychosis. *Early Interv Psychiatry*. 2015; 9(5):345–356. <https://doi.org/10.1111/eip.12193> PMID: 25345316
13. Biswas S. Digital Indians: Ben Gomes. BBC news. 2013 Sept 10. Available from: <http://www.bbc.com/news/technology-23866614>. Accessed Nov 1, 2017
14. Burns JM, Davenport TA, Durkin LA, Luscombe GM, Hickie IB. The internet as a setting for mental health service utilization by young people. *Med J Aust*. 2010 Jun 7; 192(11 Suppl):S22–6. PMID: 20528703
15. Berger M, Wagner TH, Baker LC. Internet use and stigmatized illness. *Soc Sci Med*. 2005 Oct; 61(8):1821–7. Epub 2005 Apr 26. <https://doi.org/10.1016/j.socscimed.2005.03.025> PMID: 16029778
16. Birnbaum ML, Rizvi AF, Faber K, Addington J, Correll CU, Gerber C, et al. Digital Trajectories to Care in First-Episode Psychosis. *Psychiatr Serv*. 2018 Dec 1; appips201800180. <https://doi.org/10.1176/appi.ps.201800180> Epub 2018 Sep 26. PMID: 30256181
17. Birnbaum ML, Rizvi AF, Confino J, Correll CU, Kane JM. Role of social media and the Internet in pathways to care for adolescents and young adults with psychotic disorders and non-psychotic mood

- disorders. *Early Interv Psychiatry*. 2017 Aug; 11(4):290–295. Erratum in: *Early Interv Psychiatry*. 2017 Dec; 11(6):539. Epub 2015 Mar 23. <https://doi.org/10.1111/eip.12237> PMID: 25808317
18. Van Meter AR, Birnbaum ML, Rizvi A, Kane JM. Online help-seeking prior to diagnosis: Can web-based resources reduce the duration of untreated mood disorders in young people? *J Affect Disord*. 2019 Jun 1; 252: 130–4. <https://doi.org/10.1016/j.jad.2019.04.019> PMID: 30981056
  19. Paparrizos J, White RW, Horvitz E. Screening for Pancreatic Adenocarcinoma Using Signals From Web Search Logs: Feasibility Study and Results. *J Oncol Pract*. 2016; 12(8):737–744. <https://doi.org/10.1200/JOP.2015.010504> PMID: 27271506
  20. White RW, Horvitz E. Evaluation of the Feasibility of Screening Patients for Early Signs of Lung Carcinoma in Web Search Logs. *JAMA Oncol*. 2017; 3(3):398–401. <https://doi.org/10.1001/jamaoncol.2016.4911> PMID: 27832243
  21. White RW, Doraiswamy PM, Horvitz E. Detecting neurodegenerative disorders from web search signals. *NPJ Digit Med*. 2018 Apr 23; 1:8. <https://doi.org/10.1038/s41746-018-0016-6> PMID: 31304293
  22. Eichstaedt JC, Smith RJ, Merchant RM, Ungar LH, Crutchley P, Preotiuc-Pietro D, et al. Facebook language predicts depression in medical records. *Proc Natl Acad Sci*. 2018 Oct 30; 115(44):11203–8 <https://doi.org/10.1073/pnas.1802331115> PMID: 30322910
  23. De Choudhury M, Gamon M, Counts S, Horvitz E. Predicting depression via social media. *Proc Int AAAI Conf Weblogs Soc Media*. 2013 Jun 28; 128–137
  24. Reece AG, Reagan AJ, Lix KL, Dodds PS, Danforth CM, Langer EJ. Forecasting the onset and course of mental illness with Twitter data. *Sci Rep*. 2017 Oct 11; 7(1):13006 <https://doi.org/10.1038/s41598-017-12961-9> PMID: 29021528
  25. De Choudhury M, Kiciman E, Dredze M, Coppersmith G, Kumar M. Discovering shifts to suicidal ideation from mental health content in social media. *Proc CHI Hum Factors Computing Syst In: ACM Digital Library*. 2016 May 7; 2098–2110
  26. Coppersmith G, Leary R, Crutchley P, Fine A. Natural language processing of social media as screening for suicide risk. *Biomed Inform Insights*. 2018 Aug; 10:1178222618792860. <https://doi.org/10.1177/1178222618792860> PMID: 30158822
  27. De Choudhury M, Counts S, Horvitz E. Predicting postpartum changes in emotion and behavior via social media. *Proc SIGCHI Hum Factors Computing Syst In: ACM Digital Library*. 2013 Apr 27; 3267–3276
  28. Birnbaum ML, Ernala SK, Rizvi AF, De Choudhury M, Kane JM. A collaborative approach to identifying social media markers of schizophrenia by employing machine learning and clinical appraisals. *JMIR*. 2017; 19(8):e289 <https://doi.org/10.2196/jmir.7956> PMID: 28807891
  29. Birnbaum ML, Ernala SK, Rizvi AF, Arenare E, Van Meter A, Choudhury MD, et al. Detecting Relapse in Youth with Psychotic Disorders Utilizing Patient-Generated and Patient-Contributed Digital Data from Facebook. *NPJ Schizophr*. 2019 Oct 7; 5(1):17. <https://doi.org/10.1038/s41537-019-0085-9> PMID: 31591400; PMCID: PMC6779748
  30. Kirschenbaum MA, Birnbaum ML, Rizvi A, Muscat W, Patel L, Kane JM. Google search activity in early psychosis: A qualitative analysis of internet search query content in first episode psychosis. *Early Interv Psychiatry*. 2019; Epub 2019 Oct 21. <https://doi.org/10.1111/eip.12886> PMID: 31637869
  31. Chung CK, Pennebaker JW. Linguistic inquiry and word count (LIWC): pronounced “Luke,” . . . and other useful facts. In: *Applied natural language processing: Identification, investigation and resolution*. Pennsylvania: IGI Global; 2012. p. 206–229. <https://doi.org/10.4018/978-1-60960-741-8.ch012>
  32. Tausczik YR, Pennebaker JW. The psychological meaning of words: LIWC and computerized text analysis methods. *J Lang Soc Psychol*. 2010 Mar; 29(1):24–54
  33. Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B (Methodological)*. 1995; 57(1), 289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031>
  34. Seabold S, Perktold J. Statsmodels: Econometric and statistical modeling with python. In *Proceedings of the 9th Python in Science Conference*. 2010; 57, 61.
  35. Addington J, Liu L, Buchy L, Cadenhead KS, Cannon TD, Cornblatt BA, et al. North American Prodrome Longitudinal Study (NAPLS 2): The Prodromal Symptoms. *J Nerv Ment Dis*. 2015 May; 203(5):328–35. <https://doi.org/10.1097/NMD.0000000000000290> PMID: 25919383; PubMed Central PMCID: PMC4417745
  36. Knowles L, Sharma T. Identifying vulnerability markers in prodromal patients: a step in the right direction for schizophrenia prevention. *CNS Spectr*. 2004 Aug; 9(8):595–602. <https://doi.org/10.1017/s1092852900002765> PMID: 15273652
  37. Martin DJ, Smith DJ. Is there a clinical prodrome of bipolar disorder? A review of the evidence. *Expert Rev Neurother*. 2013; 13(1):89–98. <https://doi.org/10.1586/ern.12.149> PMID: 23253393

38. Sheffield JM, Karcher NR, Barch DM. Cognitive Deficits in Psychotic Disorders: A Lifespan Perspective. *Neuropsychol Rev*. 2018; 28(4):509–533. <https://doi.org/10.1007/s11065-018-9388-2> PMID: 30343458
39. Shmukler AB, Gurovich IY, Agius M, Zaytseva Y. Long-term trajectories of cognitive deficits in schizophrenia: A critical overview. *Eur Psychiatry*. 2015; 30(8):1002–1010. <https://doi.org/10.1016/j.eurpsy.2015.08.005> PMID: 26516984
40. Baglioni C, Nanovska S, Regen W, Spiegelhalter K, Feige B, Nissen C, et al. Sleep and mental disorders: A meta-analysis of polysomnographic research. *Psychol Bull*. 2016; 142(9):969–990. <https://doi.org/10.1037/bul0000053> PMID: 27416139
41. Benson KL. Sleep in Schizophrenia: Pathology and Treatment. *Sleep Med Clin*. 2015; 10(1):49–55. <https://doi.org/10.1016/j.jsmc.2014.11.001> PMID: 26055673
42. Takaesu Y. Circadian rhythm in bipolar disorder: A review of the literature. *Psychiatry Clin Neurosci*. 2018; 72(9):673–682. <https://doi.org/10.1111/pcn.12688> PMID: 29869403
43. Buck B, Minor KS, Lysaker PH. Differential lexical correlates of social cognition and metacognition in schizophrenia; a study of spontaneously-generated life narratives. *Compr Psychiatry*. 2015 Apr 1; 58:138–45 <https://doi.org/10.1016/j.comppsy.2014.12.015> PMID: 25600423
44. Buck B, Penn DL. Lexical characteristics of emotional narratives in schizophrenia: relationships with symptoms, functioning, and social cognition. *J Nerv Ment Dis*. 2015 Sep; 203(9):702 <https://doi.org/10.1097/NMD.0000000000000354> PMID: 26252823
45. Hong K, Nenkova A, March ME, Parker AP, Verma R, Kohler CG. Lexical use in emotional autobiographical narratives of persons with schizophrenia and healthy controls. *Psych Res*. 2015 Jan 30; 225(1–2):40–9
46. Minor KS, Bonfils KA, Luther L, Firmin RL, Kukla M, MacLain VR, et al. Lexical analysis in schizophrenia: how emotion and social word use informs our understanding of clinical presentation. *Psych Res*. 2015 May 1; 64:74–8
47. Fineberg SK, Leavitt J, Deutsch-Link S, Dealy S, Landry CD, Pirruccio K, et al. Self-reference in psychosis and depression: a language marker of illness. *Psychol Med*. 2016 Sep; 46(12):2605–15 <https://doi.org/10.1017/S0033291716001215> PMID: 27353541
48. Bedi G, Carrillo F, Cecchi GA, Slezak DF, Sigman M, Mota NB, et al. Automated analysis of free speech predicts psychosis onset in high-risk youths', *NPJ Schizophr*, 1, 15030. <https://doi.org/10.1038/npschz.2015.30> PMID: 27336038
49. Corcoran CM, Carrillo F, Fernández-Slezak D, Bedi G, Klim C, Javitt DC, et al. Prediction of psychosis across protocols and risk cohorts using automated language analysis. *World J Psychiatry*. 2018 Feb; 17(1):67–75
50. Rezaei N, Walker E, Wolff P. A machine learning approach to predicting psychosis using semantic density and latent content analysis. *NPJ schizophrenia*. 2019; 5. <https://doi.org/10.1038/s41537-019-0077-9> PMID: 31197184
51. Strous RD, Koppel M, Fine J, Nachliel S, Shaked G, Zivotofsky AZ. Automated characterization and identification of schizophrenia in writing. *J Nerv Ment Dis*. 2009 Aug 1; 197(8):585–8 <https://doi.org/10.1097/NMD.0b013e3181b09068> PMID: 19684495
52. de Boer JN, Voppel AE, Begemann MJ, Schnack HG, Wijnen F, Sommer IE. Clinical use of semantic space models in psychiatry and neurology: A systematic review and meta-analysis. *Neurosci Biobehav Rev*. 2018 Oct 1; 93:85–92 <https://doi.org/10.1016/j.neubiorev.2018.06.008> PMID: 29890179
53. Elvevåg B, Foltz PW, Weinberger DR, Goldberg TE. Quantifying incoherence in speech: an automated methodology and novel application to schizophrenia. *Schizophr Res*. 2007 Jul 1; 93(1–3):304–16 <https://doi.org/10.1016/j.schres.2007.03.001> PMID: 17433866
54. Elvevåg B, Foltz PW, Rosenstein M, DeLisi LE. An automated method to analyze language use in patients with schizophrenia and their first-degree relatives. *J Neurolinguistics*. 2010 May 1; 23(3):270–84 <https://doi.org/10.1016/j.jneuroling.2009.05.002> PMID: 20383310
55. Pauselli L, Halpern B, Cleary SD, Ku B, Covington MA, Compton MT. Computational linguistic analysis applied to a semantic fluency task to measure derailment and tangentiality in schizophrenia. *Psychiatry Res*. 2018 May 1; 263:74–9 <https://doi.org/10.1016/j.psychres.2018.02.037> PMID: 29502041
56. Gupta T, Hespos SJ, Horton WS, Mittal VA. Automated analysis of written narratives reveals abnormalities in referential cohesion in youth at ultra high risk for psychosis. *Schizophr Res*. 2018 Feb 1; 192:82–8 <https://doi.org/10.1016/j.schres.2017.04.025> PMID: 28454920
57. Mota NB, Vasconcelos NA, Lemos N, Pieretti AC, Kinouchi O, Cecchi GA, et al. Speech graphs provide a quantitative measure of thought disorder in psychosis. *PloS One*. 2012 Apr 9; 7(4):e34928 <https://doi.org/10.1371/journal.pone.0034928> PMID: 22506057



58. Mota NB, Copelli M, Ribeiro S. Thought disorder measured as random speech structure classifies negative symptoms and schizophrenia diagnosis 6 months in advance. *NPJ Schizophr.* 2017 Apr 13; 3(1):18. <https://doi.org/10.1038/s41537-017-0019-3> PMID: 28560264
59. Bauer M, Glenn T, Monteith S, Bauer R, Whybrow PC, Geddes J. Ethical perspectives on recommending digital technology for patients with mental illness. *Int J Bipolar Disord.* 2017 Dec; 5(1):6. <https://doi.org/10.1186/s40345-017-0073-9> PMID: 28155206