



## King's Research Portal

DOI:

<https://doi.org/10.3389/fpsyt.2020.553463>

[Link to publication record in King's Research Portal](#)

*Citation for published version (APA):*

Bittar, A., Velupillai, S., Downs, J., Sedgwick, R., & Dutta, R. (2020). Reviewing a decade of research into suicide and related behaviour using the South London and Maudsley NHS Foundation Trust Clinical Record Interactive Search (CRIS) system. *Frontiers in psychiatry / Frontiers Research Foundation*. <https://doi.org/10.3389/fpsyt.2020.553463>

### **Citing this paper**

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

### **General rights**

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

### **Take down policy**

If you believe that this document breaches copyright please contact [librarypure@kcl.ac.uk](mailto:librarypure@kcl.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.

1 **Reviewing a decade of research into suicide and related behaviour using**  
2 **the South London and Maudsley NHS Foundation Trust Clinical Record**  
3 **Interactive Search (CRIS) system**  
4

5 **André Bittar<sup>1\*</sup>, Sumithra Velupillai<sup>1</sup>, Johnny Downs<sup>1,2</sup>, Rosemary Sedgwick<sup>1,2</sup>, Rina**  
6 **Dutta<sup>1,2</sup>**

7  
8 <sup>1</sup> Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, UK  
9 <sup>2</sup> South London and Maudsley NHS Foundation Trust, London, UK  
10

11 **\* Correspondence:**  
12 Corresponding Author  
13 [andre.bittar@kcl.ac.uk](mailto:andre.bittar@kcl.ac.uk)

14 **Keywords: Electronic Health Records; Natural Language Processing; Machine**  
15 **Learning; Suicide, attempted; Suicide, completed; Self-Injurious Behavior**

16  
17 **Abstract**  
18

19 Suicide is a serious public health issue worldwide, yet current clinical methods for assessing  
20 a person's risk of taking their own life remain unreliable and new methods for assessing  
21 suicide risk are being explored. The widespread adoption of electronic health records (EHRs)  
22 has opened up new possibilities for epidemiological studies of suicide and related behaviour  
23 amongst those receiving healthcare. These types of records capture valuable information  
24 entered by healthcare practitioners at the point of care. However, much recent work has relied  
25 heavily on the structured data of EHRs, whilst much of the important information about a  
26 patient's care pathway is recorded in the unstructured text of clinical notes.  
27

28 Accessing and structuring text data for use in clinical research, and particularly for suicide  
29 and self-harm research, is a significant challenge that is increasingly being addressed using  
30 methods from the fields of natural language processing (NLP) and machine learning (ML). In  
31 this review, we provide an overview of the range of suicide-related studies that have been  
32 carried out using the Clinical Records Interactive Search (CRIS): a database for  
33 epidemiological and clinical research that contains de-identified EHRs from the South  
34 London and Maudsley NHS Foundation Trust. We highlight the variety of clinical research  
35 questions, cohorts and techniques that have been explored for suicide and related behaviour  
36 research using CRIS, including the development of NLP and ML approaches.  
37

38 We demonstrate how EHR data provides comprehensive material to study prevalence of  
39 suicide and self-harm in clinical populations. Structured data alone is insufficient and NLP  
40 methods are needed to more accurately identify relevant information from EHR data. We also  
41 show how the text in clinical notes provide signals for ML approaches to suicide risk  
42 assessment. Further studies in generalisability and external validation of these NLP and ML  
43 approaches are needed.  
44

45 **1. Introduction**

46

47 **1.1 Suicidality research prior to CRIS**

48

49 Prior to the introduction of electronic health records (EHRs), the study of suicidality in  
50 Camberwell, the southeast London catchment area served by King’s College Hospital, was  
51 undertaken by paper case note review, for example of all referrals to a self-harm team over a  
52 6 month period (Neeleman et al., 1996). Data was painstakingly extracted and checked from  
53 each consecutive referral to ensure they fitted written criteria and in the Neeleman et al.  
54 (1996) study a single research question about ethnic differences was posed.

55

56 Later, when Dutta et al. (2010) were trying to determine the epidemiology of completed  
57 suicides in a clinically representative cohort of patients experiencing their first episode of  
58 psychosis over a 40-year inception period, it was imperative that diagnostic consistency was  
59 stringent. They achieved this by amalgamating the Camberwell Cumulative Psychiatric Case  
60 Register for the period between January 1, 1965, and December 31, 1983 (Castle et al.,  
61 1991), and then for the period between January 1, 1984, and December 31, 2004, using the  
62 basic hospital computer records held at the time with structured fields, to generate a list of all  
63 patients admitted with any possible psychotic illness (according to ICD-9 and ICD-10 codes).  
64 They then used the information gleaned from reading through the paper case records of all  
65 these patients, including medical, nursing, social work, and occupational therapy notes,  
66 together with all correspondence relating to the year after each patient’s first presentation to  
67 complete the Operational Checklist for Psychotic Disorders (OPCRIT) (McGuffin et al.,  
68 1991). This is a well-validated symptom checklist which enabled operational research  
69 diagnostic criteria (RDC) (Spitzer et al., 1978) computer diagnoses to be made using the  
70 OPRCIT program.

71

72 This methodology meant inclusion in the cohort was clearly and consistently defined, and the  
73 outcome of deaths by suicide and open verdicts up until March 31, 2007 according to the  
74 International Classification of Diseases (ICD) was identified by a direct case-tracing  
75 procedure with the Office for National Statistics (ONS) for England and Wales and the  
76 General Register Office (GRO) for Scotland. OPCRIT+ (a redesigned version of OPCRIT for  
77 use in clinical settings with an expanded number of objectively rated items) facilitated access  
78 to structured symptom information entered by clinicians to generate diagnoses including  
79 ‘suicidal ideation’ but not self-harm (Rucker et al., 2011), limiting its application for the  
80 study of self-harm and suicidal behaviour.

81

82

83 **1.2 Why EHRs and CRIS?**

84

85 The widespread adoption of EHRs has meant that large-scale clinical data are now available  
86 for clinical research, although researchers have to contend with the large volume, complexity  
87 and heterogeneity of these ‘big data’ resources. Typical EHR systems store patient data in  
88 both structured fields and as unstructured text (as well as other media types, such as medical  
89 images). Structured data fields, such as drop-down menus, forms and checkboxes, tend to be  
90 made available to clinical practitioners as a means to directly encode patient diagnoses,  
91 assessment results, etc. in a predetermined format. Unstructured text entry allows for more  
92 nuanced documentation, providing context to assessments, patient status, and other  
93 information pertinent to the clinical interaction. The availability of these electronic health  
94 data has greatly facilitated mental health research. Investigators can now use EHRs to gather

95 data about clinical populations, identify participants for clinical trials, carry out retrospective  
96 case-control studies, develop and trial predictive models and guide the implementation of  
97 evidence-based practices (Casey et al., 2016; Castillo et al., 2015).  
98

99 In 2008, the South London and Maudsley National Health Service (NHS) Foundation Trust  
100 Biomedical Research Centre (SLaM BRC) developed the Clinical Record Interactive Search  
101 (CRIS) application. CRIS is a repository of anonymised, structured and free-text data derived  
102 from the electronic health record system used by SLaM [See Perera et al. (2016) for further  
103 details]. CRIS provides unprecedented information on mental disorders and outcomes in  
104 routine clinical care at scale, particularly through enhancements from the use of natural  
105 language processing (NLP) to extract previously inaccessible information, ranging from  
106 patients' cognitive function, smoking status and education, to antipsychotic medication  
107 profiles and substance misuse (Jackson et al., 2017), as well as linkages to external data  
108 sources such as national mortality data from the ONS (Hayes et al., 2014), education data  
109 (National Pupil Database) (J. Downs et al., 2019), and Hospital Episode Statistics (HES)  
110 (Chang et al., 2017). CRIS has also allowed smaller-scale linkages, such as SHIELD, a  
111 service improvement project investigating self-harm at the emergency departments of two  
112 major London hospitals (Polling et al., 2015).  
113

114 The availability of this type of large-scale data heralds the prospect of using statistical and  
115 data science approaches to analyse larger cohorts and better understand how these behaviours  
116 manifest in healthcare settings (Velupillai, Hadlaczky, et al., 2019). However, using these  
117 data also presents major challenges, as much of the key clinical information, including  
118 suicidal behaviour, is recorded as unstructured clinical case notes and correspondence  
119 (Anderson et al., 2015; Haerian et al., 2012; Metzger et al., 2017).  
120

121 Over the last 10 years, researchers have used CRIS to conduct a number of epidemiological  
122 studies to examine suicidal behaviours across a range of mental health conditions (e.g.  
123 autism, psychotic disorders), and demographic groups (e.g. adults, children and adolescents,  
124 pregnant women). Methodologies have evolved, improving the accuracy of identifying  
125 suicidality-related constructs and predictive models of suicide risk. In the following sections,  
126 we review the evidence generated from CRIS on suicidal behaviours, the NLP methods used,  
127 and the value of the resulting cohorts and datasets created.  
128

## 129 **2. Identification / prevalence estimates of suicidality in clinical populations** 130

131 Suicide-related behaviour is the manifestation of a complex set of phenomena that depend on  
132 many contextual factors which can change quickly from one day to another. Completed  
133 suicide remains relatively rare, meaning that tools to assess suicide risk must have a high  
134 predictive validity to be of use in a clinical setting (Carter et al., 2017). Accurate  
135 identification of suicide-related behaviour is, therefore, both highly challenging and of prime  
136 importance in determining prevalence of suicidal behaviour in clinical populations and for the  
137 development of risk models. While the earliest studies on suicide and related behaviour in  
138 CRIS relied on structured fields and mortality data linkages to identify cohorts, increasing  
139 efforts have focussed on using NLP to identify suicidality-related concepts in the high  
140 volume of unstructured clinical text held in the database. The task of automatically  
141 identifying mentions of suicidal behaviour in clinical notes is complicated by the necessity to  
142 distinguish actual events relating to the patient from negated mentions, behaviour reported as  
143 family history, or those that are recorded with a degree of uncertainty (Velupillai et al.,  
144 2018). Furthermore, given the inherent variation across clinical populations, which is

145 reflected in the language used in clinical reporting, NLP tools developed for one population  
146 may not be reliably transferable to another without adaptation. NLP systems used to identify  
147 suicide-related constructs in clinical notes must, therefore, be developed for and validated  
148 within each target population.

149  
150 A wide range of known risk and contributory factors are associated with suicide, with  
151 symptoms of mental illness being recognisable in more than 85% of people who die by  
152 suicide, according to psychological autopsy interviews with family, friends and medical  
153 professionals (Arsenault-Lapierre et al., 2004; Cavanagh et al., 2003). Over the last ten years,  
154 research using CRIS has been conducted to examine the associations of self-harm, suicidality  
155 and death by suicide with mental health conditions and a broad range of situational factors,  
156 from homelessness to drug misuse to poor continuity of service (Tulloch et al., 2012;  
157 Bogdanowicz et al., 2016; Lopez-Morinigo et al., 2014). As we describe in our summary  
158 below, initial studies on suicide and related behaviour in CRIS used structured fields held  
159 within standard assessment forms or diagnostic codes. Progressively, researchers began to  
160 make use of CRIS's free-text fields and search functionalities, while more recently, NLP  
161 techniques have been employed to extract and structure suicide-related information from  
162 within the case notes. The principal characteristics of the clinical cohorts mentioned in this  
163 review are summarised in [Table 1](#).

Deleted: Table 1

## 165 **2.1 Using structured data**

166

### 167 **2.1.1 Suicidality outcome data**

168

169 The Health of the Nation Outcome Scales (HoNOS) were introduced in 1996, to measure the  
170 health and social functioning of people with mental illness. Within SLAM, as with most UK  
171 mental health trusts, clinicians are expected to complete HoNOS for all patients receiving  
172 care. The non-accidental self-injury item on the HoNOS score has been shown to be the only  
173 individual item associated with higher mental health service costs (Twomey et al., 2016). It  
174 has been used in a number of studies in CRIS to assess both the direct and indirect impact of  
175 self-harm. The individual non-accidental self-injury HoNOS item has been included as a  
176 covariate in a number of analyses of adverse outcomes within CRIS. These include  
177 homelessness and length of hospital stay for psychiatric inpatients (Tulloch et al., 2012),  
178 functional status and mortality in serious mental illness (Hayes et al., 2012), facilitated  
179 discharge and bed days (A. D. Tulloch et al., 2015), and the effects of clozapine on premature  
180 mortality (Hayes et al., 2014). When assessing self-harm as a potential risk factor for  
181 mortality among patients with personality disorder, the HoNOS item was again used in  
182 isolation as a marker of self-harm risk (Fok et al., 2014). Despite the provision of optional  
183 structured questionnaires on CRIS, such as the Patient Health Questionnaire-9 (PHQ-9)  
184 (whose final item enquires about thoughts of self-harm and suicide) and the Beck Scale for  
185 Suicide Ideation (BSS), very few are completed in general clinical work where free-text input  
186 is favoured by clinicians, making them of limited value for studies of real-world clinical  
187 cohorts. Conversely local NHS Trust requirements to complete structured suicide risk  
188 assessments for all patients means this data is better recorded and has been studied.

189

### 190 **2.1.2 Suicide risk assessment data**

191

192 Structured suicide and violence risk assessments in mental health services has been shown to  
193 have low predictive accuracy for all-cause mortality (Lopez-Morinigo et al., 2016), however  
194 these assessments have continued to be used in clinical practice. Lopez-Morinigo et al.

196 examined the use of risk assessment proforma for their investigation into suicide completion  
197 in secondary mental health care. The risk proforma, which clinicians were expected to use at  
198 that time according to local clinical policy, consisted of present/absent tick boxes for factors  
199 including suicidal history, suicidal ideation and alcohol misuse. They found that patients with  
200 a diagnosis other than schizophrenia spectrum disorder who had died by suicide, were much  
201 less likely than patients with schizophrenia to either have had a full risk assessment or a  
202 complete HoNOS even though they showed increased frequency and greater predictability in  
203 key suicide risk assessment factors: suicidal ideation, hopelessness, impulsivity and  
204 significant loss (Lopez-Morinigo et al., 2014). In their later study, they found structured risk  
205 assessment relating to suicide in schizophrenia spectrum disorders to be of little use in  
206 predicting completed suicide, with risk assessments fully completed in only 43.6% of patients  
207 who had died by suicide (Lopez-Morinigo et al., 2016). Subsequent work revealed a limited  
208 role for structured risk assessment, especially in its usefulness in revealing more nuanced  
209 factors relevant to suicide risk such as ‘mental pain’ (Lopez-Morinigo et al., 2018). They  
210 suggest that research should “switch the focus from long-term risk factors to short-term risk  
211 algorithms, which are more relevant to the clinician”.

212

### 213 **2.1.3 Suicide mortality data**

214

215 Research into mortality, including death by suicide, has typically utilised ICD-10 diagnostic  
216 codes (which must be completed as part of clinical) assessment, linked with outcome data  
217 from the Office of National Statistics (Das-Munshi et al., 2017; Hayes et al., 2014). In a  
218 retrospective cohort study, Roberts et al. (2016) used CRIS to investigate the mortality of  
219 individuals in secondary and tertiary care who had been diagnosed with CFS. Although all-  
220 cause mortality for people with CFS was not significantly different to that of the general  
221 population, there was a significantly elevated risk of completed suicide. CRIS has also been  
222 used to conduct a number of pharmaco-epidemiological studies, for example Hayes et al.  
223 (2014) examined the risk or potential risk mitigation of psychopharmacological interventions  
224 on death by suicide. Findings demonstrated treatment with the medication clozapine was  
225 associated with a reduction in risk of death by unnatural causes, including suicide, as well as  
226 natural causes.

227

## 228 **2.2 Using unstructured data**

229

### 230 **2.2.1 Free-text keywords to study self-harm presentations to emergency departments**

231

232 Polling et al. (2015) used external data linkages in combination with CRIS data (including  
233 keywords recorded in free-text fields) to create a novel dataset for the study of self-harm,  
234 which is strongly associated with mental health disorders and is the strongest single risk  
235 factor for future suicide. In England, population-level assessment of self-harm is recorded in  
236 the Hospital Episode Statistics (HES) database. However, many emergency department  
237 attendances, namely those that do not lead to a hospital admission, still go unrecorded in  
238 HES, and completion of the reason for presentation is low, thus limiting the value of this data  
239 source for studies of self-harm presentations. Polling et al. addressed these shortcomings by  
240 combining routinely collected data from electronic health records in CRIS and HES. They  
241 validated their data against another dataset curated through manual review of emergency  
242 department notes and audit forms, also compiling a list of self-harm search terms.

243

244 **2.2.2 Free-text keywords to study perinatal self-harm in women with psychiatric**  
245 **disorders**

246  
247 Using the self-harm-related terms identified by Polling et al. (2015), Taylor et al. (2016)  
248 investigated the prevalence and risk factors of self-harm and suicide ideation in women with  
249 psychotic disorders and bipolar disorder during pregnancy. They identified a cohort of 420  
250 patients by performing a free-text search of CRIS records for both suicidal ideation and self-  
251 harm. The perinatal period is generally associated with lower risk of both suicide and self-  
252 harm in the general population, however, women diagnosed with severe postpartum  
253 psychiatric disorders are up to 70 times more at risk of suicide. In Taylor et al.'s cohort,  
254 24.3% of women had a report of suicidal ideation and 7.9% had a recorded self-harm event  
255 during their index pregnancy.

256  
257 **2.2.3 Free-text keywords to study self-harm and human trafficking**

258  
259 In a further study using the free-text search capabilities of CRIS, Borschmann et al. (2017)  
260 carried out an analysis of self-harm among victims of human trafficking. They identified  
261 patients for their cohort by searching the CRIS free-text notes for terms indicating possible  
262 trafficking (e.g. "victim of trafficking", "sex trafficking", "trafficked"). In the same way,  
263 documents were screened for mentions of self-harm behaviour using a list of terms including  
264 "self-harm", "DSH", "burn\*" and "electrocute\*". They found that 33% of all trafficked  
265 patients had engaged in self-harm prior to care, while 25% did so during care. After self-  
266 harming, trafficked patients were subsequently more likely to be admitted to a ward than  
267 those who had not been victims of human trafficking. After self-harming, trafficked patients  
268 were more likely than non-trafficked patients to be admitted as a psychiatric inpatient, but  
269 less likely to attend an emergency department.

270  
271 **2.3 Using Natural Language Processing (NLP)**

272  
273 **2.3.1 Study of mortality in opioid use disorder patients using NLP to identify cohorts**

274  
275 Using data from CRIS with an external linkage to ONS mortality data, Bogdanowicz et al.  
276 (2016) investigated the effectiveness of addiction-specific clinical risk assessments for  
277 identifying groups with high mortality in opioid use disorder (OUD). Patients with a  
278 diagnosis of OUD were identified by ICD-10 code F11. ICD-10 diagnosis was supplemented  
279 with structured output of one of the CRIS NLP tools that identifies diagnoses in unstructured  
280 clinical notes. Overdose (both accidental and intentional) was the most common cause of  
281 death and clinically assessed suicidality was found to be significantly associated with  
282 increased overdose mortality.

283  
284 **2.3.2 NLP to identify suicide-related behaviour**

285  
286 Today, with the increasing body of research on suicide and related behaviour in CRIS, and a  
287 diversity of clinical population groups under study, has come a need to develop more targeted  
288 methods of accessing the suicide-related data within the unstructured clinical narratives. NLP  
289 systems designed for this task need to identify the different types of suicide-related behaviour  
290 (suicide attempt, suicidal ideation, self-harm, etc.) and account for the linguistic variation that  
291 indicates whether a mention is attested, negated or uncertain, is relevant to the patient, or a  
292 family member, and so on. These considerations have spurred on the recent development of  
293 bespoke NLP tools. For example, Gkotsis et al. (2016) developed an NLP system specifically

294 designed to detect whether a suicide-related concept is negated or not. This system was  
295 developed and evaluated on a random sample of clinical notes from CRIS. In a more recent  
296 study, Fernandes et al. (2018) developed two NLP approaches to detect relevant mentions of  
297 suicidal ideation and another to identify recorded suicide attempts.  
298

### 299 **2.3.3 NLP features to identify key suicide risk periods**

300 Identifying periods during which a patient is at elevated risk of making a suicide attempt is  
301 key to enabling timely intervention. However, information available to clinicians concerning  
302 the rapidly changing dynamic factors leading up to a suicide attempt has been limited. Bittar  
303 et al. (2019) explored whether it is possible to use EHRs to automatically predict suicide  
304 attempts in a broad clinical population (across all age groups) using only data from a  
305 relatively short period of 30 days leading up to an event. This work was based on the  
306 hypothesis that periods prior to a suicide attempt are a time of acute crisis that is reflected,  
307 explicitly or implicitly, in clinician records, making these periods distinguishable from  
308 periods not preceding an attempted suicide. Combining all three features of (1) structured  
309 data from EHRs, (2) structured values extracted by NLP software, and (3) vectorised bag-of-  
310 words of all documents provided the best model to classify or distinguish between “document  
311 windows” prior to a suicide attempt or not.  
312

### 314 **2.3.4 NLP to study suicidal behaviour in children and adolescents**

315  
316 The risk and conceptualisation of suicidal behaviour for children and adolescents can be  
317 different to adults (Cha et al., 2018). Downs et al. (2017) conclude that the clinician notes on  
318 suicidal risk in children and adolescents are different to an adult review. For example,  
319 clinicians may have a greater reliance on third person report, where caregivers voice concerns  
320 regarding the young person’s suicidality. It is also possible that suicidality is ‘discovered’  
321 rather than being the presenting complaint, hence changing the emphasis and position of  
322 suicide-related text/progress notes within the young person’s clinical record.  
323

324 Adolescence is associated with a high risk of suicide and self-harm compared to most other  
325 age groups, but few studies have examined the prevalence of suicidal behaviour in large  
326 adolescent patient cohorts. Downs et al. (2017) first used CRIS to explore suicidality in  
327 young people but focused on a population with autism spectrum disorders (ASD), who have  
328 shown much greater risk of suicidal behaviours than neurotypically developing children. A  
329 cohort of young people diagnosed with ASD were identified and NLP techniques were used  
330 to identify suicidal behaviour from the clinical notes in CRIS. Their corresponding free-text  
331 notes (progress reports, medical correspondence, risk assessments, etc.) were manually  
332 annotated for mentions of suicidality by clinical researchers. A prevalence analysis of  
333 suicidality in a sample of the data showed that only 3% of all documents mentioned suicide-  
334 related information.  
335

336 Using a subset of this cohort, Holden et al. (2020) used a historical cohort design and applied  
337 NLP approaches to extract information on victimisation by bullying and suicidal behaviour.  
338 They found those young people with ASD who were bullied were nearly twice as likely to  
339 report later suicidal ideation. The dataset created by Downs et al. has also recently proven  
340 useful to train machine learning models for use in suicide research. Song et al. (2020) used a  
341 revised version of the data to develop a deep neural network classifier that identifies  
342 sentences containing positive mentions of suicidality while taking into account the contextual  
343 information in surrounding sentences. This type of approach provides an alternative to



344 modelling suicide-related information from text that better takes into account the narrative  
345 discourse in the clinical documentation.

346

347 Velupillai, Epstein, et al. (2019) developed and validated a method for identifying suicidality  
348 across a more heterogenous clinical population in EHRs using NLP, expanding the  
349 population beyond ASD. They examined 1,601,422 documents from 23,455 young people  
350 and developed a method to accurately identify suicidal behaviour information in a very broad  
351 clinical population. The resulting dataset and NLP approaches used, provide a powerful  
352 example of how NLP approaches can be used to rapidly examine the prevalence of suicidal  
353 behaviour in very large adolescent clinical populations.

354

### 355 **2.3.5 NLP to study depression and suicidality in older adults**

356

357 Free-text mentions of depressive symptoms were used as outcome measures in the  
358 assessment of later-life depression in people from ethnic minorities by Mansour et al. (2020).  
359 This study used NLP tools designed to detect depressive symptoms recorded in unstructured  
360 texts in CRIS, including the identification of mentions of suicidal ideation. These depressive  
361 symptom NLP tools, developed to account for the presence of contextual markers such as  
362 negation and irrelevant concepts, were also used by Cai et al. (2020) in their investigation  
363 into predictors of mortality in people with late life depression.

364

## 365 **3. The Next Ten Years?**

366

367 Although EHR data are not created for research purposes, they provide a rich resource for  
368 large-scale retrospective research, allowing identification of diverse and comprehensive  
369 clinical study samples. One of the main challenges in suicide research is obtaining  
370 sufficiently large study samples to study an outcome with a high enough base rate for  
371 predictive modelling to have a meaningful positive predictive value. The low base rate of  
372 completed suicide limits the predictive value of any model (whether established statistical  
373 techniques or machine learning) (McHugh & Large, 2020), but related behaviours, such as  
374 suicidal ideation, intention, planning and self-harm can be studied. Over the past 10 years,  
375 CRIS has provided an unprecedented resource for studying suicide and related behaviour in a  
376 UK clinical population to an extent that would not have been possible before the introduction  
377 of EHRs.

378

379 Furthermore, EHRs reflect real-world clinical practice. This means that the context of how,  
380 for example, structured risk assessment tools and other schedules, like HoNOS, are used in  
381 daily clinical work needs to be well understood when including them as variables in clinical  
382 research studies. Most of the relevant information is found in the free text, and appropriate  
383 NLP solutions are key components for enabling risk modelling.

384

385 Looking to the future, validation of the various findings, including the developed NLP  
386 applications, on other EHR systems and in other clinical catchment areas would provide  
387 insights into the generalisability of these models to new clinical settings. CRIS has also been  
388 implemented in other sites, such as the Camden & Islington Research Database (Werbeloff et  
389 al., 2018). Comparisons of algorithms' performance in such datasets would further advance  
390 this field and provide evidence about the broader generalisability of findings. Collaborative  
391 efforts are currently being made to compare methodologies and NLP tools across healthcare  
392 institutions not just within the UK, but also with collaborators in the USA.

393

394 Furthermore, advances in computational analysis of EHR data, e.g. machine learning in  
395 combination with NLP, will continue to develop, and provide novel solutions to suicide  
396 research (Walsh et al., 2017). With the existing CRIS subsets, clinical cohorts, and NLP  
397 approaches developed for the studies described in this review, benchmarks have been created  
398 that allow for appropriate comparisons between different methodologies.  
399

400 Going beyond identification or prediction of those at risk, analysis of continuously collected  
401 data, and integration of EHR data with smartphone, wearable device and even social media  
402 data could allow collection of data across different time periods, not just at the time of  
403 clinical interactions, thus helping to understand suicidal crises and enabling delivery of  
404 targeted suicide prevention interventions (Torous & Walker, 2019).  
405

#### 406 **4. Conclusion**

407  
408 In this review of a decade of research into suicide and related behaviour using CRIS we have  
409 summarised the evolution of different methods employed to identify suicide and related  
410 behaviour, including linkages to mortality data, structured ICD-10 codes, manual review of  
411 clinical notes, keyword searching in free text and relevant mentions identified using NLP  
412 techniques. Cohorts under study have varied in size from several hundred to tens of  
413 thousands of patients and have covered adult, elderly as well as child and adolescent patients.  
414 A range of clinical disorders have been described from the perspective of suicide and related  
415 behaviours, including pregnancy, severe mental illness and self-harm, opioid use disorder  
416 patients, chronic fatigue syndrome and autism spectrum disorders. Finally, some studies have  
417 identified and investigated specific clinical events, such as emergency department  
418 attendances or hospital admissions.  
419

420 Research in NLP methods has evolved over the years, from methods relying on symbolic  
421 principles (e.g. grammars, lexicon development, pattern matching) to statistical methods and,  
422 more recently, machine learning approaches. This progression is also reflected in the NLP  
423 work that has been developed using CRIS.

424 The first approaches that were developed to process CRIS data were pattern matching  
425 approaches to identify certain pieces of information (e.g. medication, smoking status,  
426 substance misuse) using the GATE framework (Cunningham et al., 2013). In many cases, the  
427 information of interest is a particular clinical construct (e.g. hallucinations, echolalia) or a  
428 specific diagnosis. A bespoke application, called TextHunter (Jackson et al., 2017), was  
429 developed for these types of constructs. TextHunter is a software application that requires a  
430 set of manually pre-annotated examples to train a supervised machine learning classifier  
431 (Support Vector Machine). These NLP applications identify and classify the relevant  
432 constructs and produce structured variables that are then stored in table columns in CRIS.  
433 Researchers may access these variables (along with the “standard” structured fields from the  
434 EHR) through the SQL interface of the CRIS database to identify cohorts of patients for  
435 epidemiological studies and clinical research. Several studies cited herein have made use of  
436 these structured variables (Bittar et al., 2019; Bogdanowicz et al., 2016; Roberts et al., 2016).  
437

438 In addition to these “integrated” NLP applications, clinicians have worked alongside NLP  
439 researchers to develop custom NLP tools to identify suicide-related constructs in specific  
440 population samples within CRIS. As we have seen, the focus of most work has been the  
441 epidemiology and prevalence of suicidal behaviour, with NLP tools that use both rule-based  
442 (Gkotsis et al., 2016; Velupillai, Epstein, et al., 2019) and machine learning paradigms  
443 (Fernandes et al., 2018), including neural network architectures (Song et al., 2020). Most

444 recently efforts have also been made to model dynamic suicide risk using supervised machine  
445 learning (Bittar et al., 2019).

446

#### 447 **Gaining Access to CRIS**

448

449 The de-identified CRIS database has received ethical approval for secondary analysis: Oxford  
450 REC C, reference 18/SC/0372. The data is used in an anonymised and data-secure format  
451 under strict governance procedures. CRIS data is made available to researchers with  
452 appropriate credentials (provided by the South London and Maudsley NHS Trust) working on  
453 approved projects. Projects are approved by a CRIS Oversight Committee, a body set up by  
454 and reporting to the South London and Maudsley Caldicott Guardian. On request, and after  
455 appropriate credentials have been obtained as well as arrangements with the lead of the  
456 respective CRIS project, data presented in this study can be viewed within the secure system  
457 firewall.

458

459 **5. References**

460

461 Anderson, H. D., Pace, W. D., Brandt, E., Nielsen, R. D., Allen, R. R., Libby, A. M., West, D. R., &  
462 Valuck, R. J. (2015). Monitoring Suicidal Patients in Primary Care Using Electronic Health Records.  
463 *The Journal of the American Board of Family Medicine*, 28(1), 65–71.

464 <https://doi.org/10.3122/jabfm.2015.01.140181>

465 Arsenaault-Lapierre, G., Kim, C., & Turecki, G. (2004). Psychiatric diagnoses in 3275 suicides: A  
466 meta-analysis. *BMC Psychiatry*, 4, 37–37. PubMed. <https://doi.org/10.1186/1471-244X-4-37>

467 Bittar, A., Velupillai, S., Roberts, A., & Dutta, R. (2019). Text Classification to Inform Suicide Risk  
468 Assessment in Electronic Health Records. *Studies in Health Technology and Informatics*, 40–44.  
469 <https://doi.org/10.3233/SHTI190179>

470 Bogdanowicz, K. M., Stewart, R., Chang, C.-K., Downs, J., Khondoker, M., Shetty, H., Strang, J., &  
471 Hayes, R. D. (2016). Identifying mortality risks in patients with opioid use disorder using brief  
472 screening assessment: Secondary mental health clinical records analysis | Elsevier Enhanced Reader.  
473 *Drug and Alcohol Dependence*, 164, 82–88. <https://doi.org/10.1016/j.drugaledep.2016.04.036>

474 Borschmann, R., Oram, S., Kinner, S. A., Dutta, R., Zimmerman, C., & Howard, L. M. (2017). Self-  
475 Harm Among Adult Victims of Human Trafficking Who Accessed Secondary Mental Health Services  
476 in England. *Psychiatric Services*, 68(2), 207–210. <https://doi.org/10.1176/appi.ps.201500509>

477 Cai, W., Mueller, C., Shetty, H., Perera, G., & Stewart, R. (2020). Predictors of mortality in people  
478 with late-life depression: A retrospective cohort study. *Journal of Affective Disorders*.  
479 <https://doi.org/10.1016/j.jad.2020.01.021>

480 Carter, G., Milner, A., McGill, K., Pirkis, J., Kapur, N., & Spittal, M. J. (2017). Predicting suicidal  
481 behaviours using clinical instruments: Systematic review and meta-analysis of positive predictive  
482 values for risk scales. *British Journal of Psychiatry*, 210(6), 387–395. Cambridge Core.  
483 <https://doi.org/10.1192/bjp.bp.116.182717>

484 Casey, J. A., Schwartz, B. S., Stewart, W. F., & Adler, N. E. (2016). Using Electronic Health Records  
485 for Population Health Research: A Review of Methods and Applications. *Annual Review of Public  
486 Health*, 37(1), 61–81. <https://doi.org/10.1146/annurev-publhealth-032315-021353>

487 Castillo, E. G., Olfson, M., Pincus, H. A., Vawdrey, D., & Stroup, T. S. (2015). Electronic Health  
488 Records in Mental Health Research: A Framework for Developing Valid Research Methods.  
489 *Psychiatric Services*, 66(2), 193–196. <https://doi.org/10.1176/appi.ps.201400200>

490 Castle, D., Wessely, S., Der, G., & Murray, R. M. (1991). The Incidence of Operationally Defined  
491 Schizophrenia in Camberwell, 1965–84. *British Journal of Psychiatry*, 159(6), 790–794. Cambridge  
492 Core. <https://doi.org/10.1192/bjp.159.6.790>

493 Cavanagh, J. T. O., Carson, A. J., Sharpe, M., & Lawrie, S. M. (2003). Psychological autopsy studies  
494 of suicide: A systematic review. *Psychological Medicine*, 33(3), 395–405. Cambridge Core.  
495 <https://doi.org/10.1017/S0033291702006943>

496 Cha, C. B., Franz, P. J., M Guzmán, E., Glenn, C. R., Kleiman, E. M., & Nock, M. K. (2018). Annual  
497 Research Review: Suicide among youth—Epidemiology, (potential) etiology, and treatment. *Journal  
498 of Child Psychology and Psychiatry, and Allied Disciplines*, 59(4), 460–482. PubMed.  
499 <https://doi.org/10.1111/jcpp.12831>

500 Chang, C.-K., Chen, C.-Y., Broadbent, M., Stewart, R., & O'Hara, J. (2017). Hospital admissions for  
501 respiratory system diseases in adults with intellectual disabilities in Southeast London: A register-  
502 based cohort study. *BMJ Open*, 7(3), e014846. <https://doi.org/10.1136/bmjopen-2016-014846>

503 Cunningham, H., Tablan, V., Roberts, A., & Bontcheva, K. (2013). Getting More Out of Biomedical  
504 Documents with GATE's Full Lifecycle Open Source Text Analytics. *PLoS Computational Biology*,  
505 9(2), e1002854. <https://doi.org/10.1371/journal.pcbi.1002854>

506 Das-Munshi, J., Chang, C.-K., Dutta, R., Morgan, C., Nazroo, J., Stewart, R., & Prince, M. J. (2017).  
507 Ethnicity and excess mortality in severe mental illness: A cohort study. *The Lancet Psychiatry*, *4*(5),  
508 389–399. [https://doi.org/10.1016/S2215-0366\(17\)30097-4](https://doi.org/10.1016/S2215-0366(17)30097-4)  
509 Downs, J., Ford, T., Stewart, R., Epstein, S., Shetty, H., Little, R., Jewell, A., Broadbent, M.,  
510 Deighton, J., Mostafa, T., Gilbert, R., Hotopf, M., & Hayes, R. (2019). An approach to linking  
511 education, social care and electronic health records for children and young people in South London: A  
512 linkage study of child and adolescent mental health service data. *BMJ Open*, *9*(1), e024355.  
513 <https://doi.org/10.1136/bmjopen-2018-024355>  
514 Downs, J., Velupillai, S., Gkotsis, G., Holden, R., Kikoler, M., Dean, H., Fernandes, A., & Dutta, R.  
515 (2017). Detection of Suicidality in Adolescents with Autism Spectrum Disorders: Developing a  
516 Natural Language Processing Approach for Use in Electronic Health Records. *Proceedings of the*  
517 *AMIA Annual Symposium*, 641–649.  
518 Dutta, R., Murray, R. M., Hotopf, M., Allardyce, J., Jones, P. B., & Boydell, J. (2010). Reassessing  
519 the Long-term Risk of Suicide After a First Episode of Psychosis. *Archives of General Psychiatry*,  
520 *67*(12), 1230–1237. <https://doi.org/10.1001/archgenpsychiatry.2010.157>  
521 Fernandes, A. C., Dutta, R., Velupillai, S., Sanyal, J., Stewart, R., & Chandran, D. (2018). Identifying  
522 Suicide Ideation and Suicidal Attempts in a Psychiatric Clinical Research Database using Natural  
523 Language Processing. *Scientific Reports*, *8*(1). <https://doi.org/10.1038/s41598-018-25773-2>  
524 Fok, M. L.-Y., Stewart, R., Hayes, R. D., & Moran, P. (2014). Predictors of Natural and Unnatural  
525 Mortality among Patients with Personality Disorder: Evidence from a Large UK Case Register. *PLOS*  
526 *ONE*, *9*(7), e100979. <https://doi.org/10.1371/journal.pone.0100979>  
527 Gkotsis, G., Velupillai, S., Oelrich, A., Dean, H., Liakata, M., & Dutta, R. (2016). Don't Let Notes  
528 Be Misunderstood: A Negation Detection Method for Assessing Risk of Suicide in Mental Health  
529 Records. *Proceedings of the Third Workshop on Computational Linguistics and Clinical Psychology*,  
530 95–105. <https://doi.org/10.18653/v1/W16-0310>  
531 Haerian, K., Salmasian, H., & Friedman, C. (2012). Methods for identifying suicide or suicidal  
532 ideation in EHRs. *AMIA ... Annual Symposium Proceedings. AMIA Symposium, 2012*, 1244–1253.  
533 PubMed.  
534 Hayes, R. D., Chang, C.-K., Fernandes, A. C., Begum, A., To, D., Broadbent, M., Hotopf, M., &  
535 Stewart, R. (2012). Functional Status and All-Cause Mortality in Serious Mental Illness. *PLOS ONE*,  
536 *7*(9), e44613. <https://doi.org/10.1371/journal.pone.0044613>  
537 Hayes, R. D., Downs, J., Chang, C.-K., Jackson, R. G., Shetty, H., Broadbent, M., Hotopf, M., &  
538 Stewart, R. (2014). The Effect of Clozapine on Premature Mortality: An Assessment of Clinical  
539 Monitoring and Other Potential Confounders. *Schizophrenia Bulletin*, *41*(3), 644–655.  
540 <https://doi.org/10.1093/schbul/sbu120>  
541 Holden, R., Mueller, J., McGowan, J., Sanyal, J., Kikoler, M., Simonoff, E., Velupillai, S., & Downs,  
542 J. (2020). Investigating Bullying as a Predictor of Suicidality in a Clinical Sample of Adolescents with  
543 Autism Spectrum Disorder. *Autism Research*, *n/a*(n/a). <https://doi.org/10.1002/aur.2292>  
544 Jackson, R. G., Patel, R., Jayatileke, N., Kolliakou, A., Ball, M., Gorrell, G., Roberts, A., Dobson, R.  
545 J., & Stewart, R. (2017). Natural language processing to extract symptoms of severe mental illness  
546 from clinical text: The Clinical Record Interactive Search Comprehensive Data Extraction (CRIS-  
547 CODE) project. *BMJ Open*, *7*(1), e012012. <https://doi.org/10.1136/bmjopen-2016-012012>  
548 Lopez-Morinigo, J.-D., Ayesa-Arriola, R., Torres-Romano, B., Fernandes, A. C., Shetty, H.,  
549 Broadbent, M., Dominguez-Ballesteros, M.-E., Stewart, R., David, A. S., & Dutta, R. (2016). Risk  
550 assessment and suicide by patients with schizophrenia in secondary mental healthcare: A case-control  
551 study. *BMJ Open*, *6*(9), e011929. <https://doi.org/10.1136/bmjopen-2016-011929>  
552 Lopez-Morinigo, J.-D., Fernandes, A. C., Chang, C.-K., Hayes, R. D., Broadbent, M., Stewart, R.,  
553 David, A. S., & Dutta, R. (2014). Suicide completion in secondary mental healthcare: A comparison

554 study between schizophrenia spectrum disorders and all other diagnoses. *BMC Psychiatry*, 14(1), 213.  
555 <https://doi.org/10.1186/s12888-014-0213-z>

556 Lopez-Morinigo, J.-D., Fernandes, A. C., Shetty, H., Ayesa-Arriola, R., Bari, A., Stewart, R., &  
557 Dutta, R. (2018). Can risk assessment predict suicide in secondary mental healthcare? Findings from  
558 the South London and Maudsley NHS Foundation Trust Biomedical Research Centre (SLaM BRC)  
559 Case Register. *Social Psychiatry and Psychiatric Epidemiology*, 53(11), 1161–1171.  
560 <https://doi.org/10.1007/s00127-018-1536-8>

561 Mansour, R., Tsamakidis, K., Rizos, E., Perera, G., Das-Munshi, J., Stewart, R., & Mueller, C. (2020).  
562 Late-life depression in people from ethnic minority backgrounds: Differences in presentation and  
563 management. *Journal of Affective Disorders*, 264, 340–347. <https://doi.org/10.1016/j.jad.2019.12.031>

564 McGuffin, P., Farmer, A., & Harvey, I. (1991). A Polydiagnostic Application of Operational Criteria  
565 in Studies of Psychotic Illness: Development and Reliability of the OPCRIT System. *Archives of*  
566 *General Psychiatry*, 48(8), 764–770. <https://doi.org/10.1001/archpsyc.1991.01810320088015>

567 McHugh, C. M., & Large, M. M. (2020). Can machine-learning methods really help predict suicide?  
568 *Current Opinion in Psychiatry, Publish Ahead of Print*. [https://journals.lww.com/co-](https://journals.lww.com/co-psychiatry/Fulltext/publishahead/Can_machine_learning_methods_really_help_predict.99128.aspx)  
569 [psychiatry/Fulltext/publishahead/Can\\_machine\\_learning\\_methods\\_really\\_help\\_predict.99128.aspx](https://journals.lww.com/co-psychiatry/Fulltext/publishahead/Can_machine_learning_methods_really_help_predict.99128.aspx)

570 Metzger, M.-H., Tvardik, N., Gicquel, Q., Bouvry, C., Poulet, E., & Potinet-Pagliaroli, V. (2017). Use  
571 of emergency department electronic medical records for automated epidemiological surveillance of  
572 suicide attempts: A French pilot study: text-mining and epidemiology of suicide attempts.  
573 *International Journal of Methods in Psychiatric Research*, 26(2), e1522.  
574 <https://doi.org/10.1002/mpr.1522>

575 Neeleman, J., Jones, P., Van Os, J., & Murray, R. M. (1996). Parasuicide in Camberwell-ethnic  
576 differences. *Social Psychiatry and Psychiatric Epidemiology*, 31(5), 284–287.  
577 <https://doi.org/10.1007/BF00787921>

578 Perera, G., Broadbent, M., Callard, F., Chang, C.-K., Downs, J., Dutta, R., Fernandes, A., Hayes, R.  
579 D., Henderson, M., Jackson, R., Jewell, A., Kadra, G., Little, R., Pritchard, M., Shetty, H., Tulloch,  
580 A., & Stewart, R. (2016). Cohort profile of the South London and Maudsley NHS Foundation Trust  
581 Biomedical Research Centre (SLaM BRC) Case Register: Current status and recent enhancement of  
582 an Electronic Mental Health Record-derived data resource. *BMJ Open*, 6(3), e008721.  
583 <https://doi.org/10.1136/bmjopen-2015-008721>

584 Polling, C., Tulloch, A., Banerjee, S., Cross, S., Dutta, R., Wood, D. M., Dargan, P. I., & Hotopf, M.  
585 (2015). Using routine clinical and administrative data to produce a dataset of attendances at  
586 Emergency Departments following self-harm. *BMC Emergency Medicine*, 15(1), 15.  
587 <https://doi.org/10.1186/s12873-015-0041-6>

588 Roberts, E., Wessely, S., Chalder, T., Chang, C.-K., & Hotopf, M. (2016). Mortality of people with  
589 chronic fatigue syndrome: A retrospective cohort study in England and Wales from the South London  
590 and Maudsley NHS Foundation Trust Biomedical Research Centre (SLaM BRC) Clinical Record  
591 Interactive Search (CRIS) Register. *The Lancet*, 387(10028), 1638–1643.  
592 [https://doi.org/10.1016/S0140-6736\(15\)01223-4](https://doi.org/10.1016/S0140-6736(15)01223-4)

593 Rucker, J., Newman, S., Gray, J., Gunasinghe, C., Broadbent, M., Brittain, P., Baggaley, M., Denis,  
594 M., Turp, J., Stewart, R., Lovestone, S., Schumann, G., Farmer, A., & McGuffin, P. (2011).  
595 OPCRIT+: An electronic system for psychiatric diagnosis and data collection in clinical and research  
596 settings. *The British Journal of Psychiatry: The Journal of Mental Science*, 199(2), 151–155.  
597 PubMed. <https://doi.org/10.1192/bjp.bp.110.082925>

598 Song, X., Downs, J., Velupillai, S., Holden, R., Kikoler, M., Bontcheva, K., Dutta, R., & Roberts, A.  
599 (2020). Using deep neural networks with intra- and inter-sentence context to classify suicidal  
600 behaviour. *Proceedings of the Twelfth International Conference on Language Resources and*  
601 *Evaluation (LREC 2020)*. LREC 2020, Marseille, France.

602 Spitzer, R. L., Endicott, J., & Robins, E. (1978). Research Diagnostic Criteria: Rationale and  
603 Reliability. *Archives of General Psychiatry*, 35(6), 773–782.  
604 <https://doi.org/10.1001/archpsyc.1978.01770300115013>

605 Taylor, C. L., van Ravesteyn, L. M., van denBerg, M. P. L., Stewart, R. J., & Howard, L. M. (2016).  
606 The prevalence and correlates of self-harm in pregnant women with psychotic disorder and bipolar  
607 disorder. *Archives of Women's Mental Health*, 19(5), 909–915. <https://doi.org/10.1007/s00737-016-0636-2>

608  
609 Torous, J., & Walker, R. (2019). Leveraging Digital Health and Machine Learning Toward Reducing  
610 Suicide—From Panacea to Practical Tool. *JAMA Psychiatry*, 76(10), 999–1000.  
611 <https://doi.org/10.1001/jamapsychiatry.2019.1231>

612 Tulloch, A. D., Khondoker, M. R., Thornicroft, G., & David, A. S. (2015). Home treatment teams and  
613 facilitated discharge from psychiatric hospital. *Epidemiology and Psychiatric Sciences*, 24(5), 402–  
614 414. Cambridge Core. <https://doi.org/10.1017/S2045796014000304>

615 Tulloch, Alex D., Khondoker, M. R., Fearon, P., & David, A. S. (2012). Associations of homelessness  
616 and residential mobility with length of stay after acute psychiatric admission. *BMC Psychiatry*, 12(1),  
617 121. <https://doi.org/10.1186/1471-244X-12-121>

618 Twomey, C., Prina, A. M., Baldwin, D. S., Das-Munshi, J., Kingdon, D., Koeser, L., Prince, M. J.,  
619 Stewart, R., Tulloch, A. D., & Cieza, A. (2016). Utility of the Health of the Nation Outcome Scales  
620 (HoNOS) in Predicting Mental Health Service Costs for Patients with Common Mental Health  
621 Problems: Historical Cohort Study. *PLOS ONE*, 11(11), e0167103.  
622 <https://doi.org/10.1371/journal.pone.0167103>

623 Velupillai, S., Epstein, S., Bittar, A., Stephenson, T., Dutta, R., & Downs, J. (2019). Identifying  
624 Suicidal Adolescents from Mental Health Records Using Natural Language Processing. *Proceedings*  
625 *of MEDINFO 2019: Health and Wellbeing e-Networks for All*, 413–417.  
626 <http://ebooks.iospress.nl/publication/52019>

627 Velupillai, S., Hadlaczky, G., Baca-Garcia, E., Gorrell, G. M., Werbeloff, N., Nguyen, D., Patel, R.,  
628 Leightley, D., Downs, J., Hotopf, M., & Dutta, R. (2019). Risk Assessment Tools and Data-Driven  
629 Approaches for Predicting and Preventing Suicidal Behavior. *Frontiers in Psychiatry*, 10.  
630 <https://doi.org/10.3389/fpsy.2019.00036>

631 Velupillai, S., Suominen, H., Liakata, M., Roberts, A., Shah, A. D., Morley, K., Osborn, D., Hayes, J.,  
632 Stewart, R., Downs, J., Chapman, W., & Dutta, R. (2018). Using clinical Natural Language  
633 Processing for health outcomes research: Overview and actionable suggestions for future advances.  
634 *Journal of Biomedical Informatics*, 88, 11–19. <https://doi.org/10.1016/j.jbi.2018.10.005>

635 Walsh, C. G., Ribeiro, J. D., & Franklin, J. C. (2017). Predicting Risk of Suicide Attempts Over Time  
636 Through Machine Learning. *Clinical Psychological Science*, 5(3), 457–469.  
637 <https://doi.org/10.1177/2167702617691560>

638 Werbeloff, N., Osborn, D. P. J., Patel, R., Taylor, M., Stewart, R., Broadbent, M., & Hayes, J. F.  
639 (2018). The Camden & Islington Research Database: Using electronic mental health records for  
640 research. *PLOS ONE*, 13(1), e0190703. <https://doi.org/10.1371/journal.pone.0190703>

641  
642

Study	Clinical group	Population size	Number of events	Age	Date range	Method to identify suicide and related behaviour
Polling et al. (2015)	Adults attending ED	7,444	10,688 ED attendances	N/A	01/04/2009-31/12/2011	ICD-10 codes X60-X84, presence of keywords related to self-harm, suicide attempts and suicidality
Bogdanowicz et al. (2016)	Patients with opioid use disorder	5,335	N/A	15-73 years mean (SD) = 37.6 (9.07) years	01/04/2008-31/03/2014	ICD-10 codes X409-X450, Y120, Y125, F119 †
Lopez-Morinigo et al. (2016)	Patients with schizophrenia spectrum disorder	426 (71 cases, 355 controls)	N/A	Mean (SD) = 44.9 (18.0) years	01/01/2007-31/12/2013	ICD-10 codes X64, X70, X71, X78, X80, X81, X84, Y10-34
Lopez-Morinigo et al. (2018)	Patients accessing secondary mental healthcare	13,758	N/A	Mean (SD) = 41.3 (12.2) for suicide, 40.6 (11.5) for no suicide	01/01/2007-01/04/2015	ICD-10 codes X64, X70, X71, X78, X80, X81, X84, Y10-34
Roberts et al. (2016)	Individuals with chronic fatigue syndrome	2,147	N/A	Mean = 39.1 years	01/01/2007-31/12/2013	ICD-10 codes X60-X84
Taylor et al. (2016)	Perinatal women with SMI	420	N/A	Mean (SD) = 31.9 (6.2) years	01/01/2007-31/12/2011	Presence of keywords [from Polling et al., 2015] related to self-harm, suicide attempts and suicidality
Downs et al. (2017)	Children and adolescents with ASD	1,906	N/A	14-18 years	01/01/2008-31/12/2013	NLP, manual classification of suicidality-related expressions
Velupillai, Epstein, et al. (2019)	Adolescents attending CAMHS	23,455	N/A	11-17 years	01/04/2009-31/03/2016	Manual annotation of suicidality-related expressions, NLP
Bittar et al. (2019)	Patients accessing secondary mental healthcare	17,640 (2,913 cases, 14,727 controls)	21,175 admissions (4,235 cases, 16,940 controls)	Mean (SD) = 33.7 (15.6) years	02/04/2006 - 31/03/2017	X6*, X7*, X80-4*, Y1*, Y2*, Y30-4*, Y87*

Table 1: Summarised characteristics of clinical cohorts created using CRIS for the study of suicide and related behaviour. †Due to indeterminacy of intent, suicide by overdose and fatal drug poisonings are grouped together.