

## SN Social Sciences

### Identifying key challenges and needs in digital mental health moderation practices supporting users exhibiting risk behaviours to develop responsible AI tools: the case study of Kooth --Manuscript Draft--

<b>Manuscript Number:</b>	SNSS-D-22-00514R1	
<b>Full Title:</b>	Identifying key challenges and needs in digital mental health moderation practices supporting users exhibiting risk behaviours to develop responsible AI tools: the case study of Kooth	
<b>Article Type:</b>	Original Article	
<b>Section/Category:</b>	Communication & Media Studies	
<b>Funding Information:</b>	UK Research and Innovation (TAS Hub)	Dr Stuart Middleton
<b>Abstract:</b>	<p>Digital platforms for mental health and wellbeing purposes have become increasingly common to help users exhibiting risk behaviours (e.g. self-harming, eating-related disorders) across all ages, opening new frontiers in supporting vulnerable users through online moderation and digital counselling. This study stems from a larger project, which explores how responsible AI solutions can up-scale existing manual moderation approaches and better target interventions for young people who ask for help or engage in risk behaviours online. This research aims to better understand the challenges and needs of moderators and digital counsellors, i.e. the 'behind the scenes'. Through this case study, the authors intend to contribute to the development of responsible AI tools that are fit for purpose and better understand the challenges. The key focus lies on Kooth.com, the UK's leading free online confidential service offering counselling and emotional wellbeing support to young people in the UK through its online web-based and pseudo-anonymous digital platform.</p>	
<b>Corresponding Author:</b>	Anita Lavorgna University of Southampton UNITED KINGDOM	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	University of Southampton	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Elena Nichele	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Elena Nichele Anita Lavorgna Stuart Middleton	
<b>Order of Authors Secondary Information:</b>		
<b>Author Comments:</b>	<p>Dear Editor,</p> <p>We are submitting a revised (and proofread by a native speaker) version of the article entitled 'Identifying key challenges and needs in digital mental health moderation practices supporting users exhibiting risk behaviours to develop responsible AI tools: the case study of Kooth' we had previously sent for consideration in SN Social Sciences.</p> <p>We thank you for your kind consideration of our paper.</p> <p>The Authors.</p>	
<b>Response to Reviewers:</b>	Please see the file attached.	

# Identifying key challenges and needs in digital mental health moderation practices supporting users exhibiting risk behaviours to develop responsible AI tools: the case study of Kooth

## Abstract

Digital platforms for mental health and wellbeing purposes have become increasingly common to help users exhibiting risk behaviours (e.g. self-harming, eating-related disorders) across all ages, opening new frontiers in supporting vulnerable users. This study stems from a larger project, which explores how responsible AI solutions can up-scale existing manual moderation approaches and better target interventions for young people who ask for help or engage in risk behaviours online. This research aims to better understand the challenges and needs of moderators and digital counsellors, i.e. the ‘behind the scenes’. Through this case study, the authors intend to contribute to the development of responsible AI tools that are fit for purpose and better understand the challenges. The key focus lies on Kooth.com, the UK’s leading free online confidential service offering counselling and emotional wellbeing support to young people in the UK through its online web-based and pseudo-anonymous digital platform.

## Keywords

Digital moderation; digital counselling; risk behaviours; responsible AI; mental health and wellbeing

## Introduction

Using social media mechanisms and other digital tools to scale up services and maximise affordances such as anonymity for mental health and wellbeing purposes is not new (e.g., McCosker, 2018). On the one hand, over the last 15 years digital platforms have offered an important mean for improving the reach and scale of mental health support, opening new frontiers (Kivitz, 2013; National Institute of Mental Health, 2017; McCosker, 2018). Among the several benefits associated with digital technologies in health settings are those associated with information access, empowerment, and the opportunity to find supportive relationships (McCosker and Darcy, 2013; Moorhead et al., 2013; Tucker and Goodings, 2017; Saha, 2020), as well as the possibility to engage with those people hard to reach and support (Tanis, 2008; Sokol and Fisher, 2016), providing additional help to individuals with urgent or special needs. More generally, the number of people seeking health-related advice (including mental health and wellbeing advice) on digital platforms is increasing, particularly because users appreciate this different style of communication, often leading to emotional care and empathy (Lederman et al., 2014). On the other hand, it has been recognised that the success of digital platforms in health settings depends on a range of support and sociocultural factors (Hansen and Aranda, 2012). These approaches alter both the ‘expert-client relationship’ (i.e., how the worker and the client interact ‘around the information sought and given’, as the level of self-disclosure

1 increases online – see Mowlabocus et al., 2015: 5) and the ‘public-professional relationship’  
2 (Kivitz, 2013), in the sense that patients and the general public are now able to remain  
3 permanently connected with health professionals and institutions (Kivitz, 2013), in a way that  
4 ‘challenges notions of expertise, whether health, biomedical or cultural, inspiring attempts to  
5 mobilise new forms of community-oriented and personalised public health intervention  
6 through digitally mediated peer practices” (McCosker, 2018:4751).  
7  
8  
9

10 In a context where mental health organisations have limited funding and hence need to  
11 carefully choose how to allocate that funding to support services, it is important to understand  
12 how to best design digital tools and to assess their effectiveness (as discussed in McCosker,  
13 2018), but also to consider the challenges and needs experienced ‘behind the scenes’ by those  
14 relying on these digital tools for their work: we believe this is a necessary step to improve the  
15 systems already available.  
16  
17  
18

19 Our study addresses this latter point by interviewing key actors (involved in moderation,  
20 counselling, emotional wellbeing support, or managerial roles), within the frame of a broader  
21 research project (the UKRI TAS Hub-funded project *SafeSpacesNLP: Behaviour classification  
22 NLP in a socio-technical AI setting for online harmful behaviours for children and young  
23 people, 2021-2022*<sup>1</sup>) that will use these insights to explore how ‘responsible AI’ solutions  
24 (Ghallab, 2019) can support up-scaling existing manual approaches and better target  
25 interventions for young people who ask for help or engage in risk behaviours online, which can  
26 have a detrimental impact on their physical (e.g., suicide, self-harming, eating-related  
27 disorders), mental (e.g., anxiety, depression, sleep disruptions, body image distortions,  
28 cyberbullying, and ‘fear of missing out’) and/or sexual health (e.g., forced marriages, sexual  
29 exploitation).  
30  
31  
32  
33  
34  
35

36 As the next section of this paper discusses in more detail, such users can be denominated  
37 ‘vulnerable publics’ (McCosker and Wilken, 2017), as they experience and share socially  
38 sensitive and emotionally charged challenges for which they seek help, through online  
39 moderation and digital counselling. Given their role, the professionals involved in online  
40 moderation and digital counselling act as frontline service workers, as they provide a blend of  
41 care support, health services and ‘feeling management’ (see Hochschild, 2003), which can be  
42 labelled as affective, emotional and immaterial labour (see McCosker and Darcy, 2013;  
43 McCosker, 2018). Moreover, since the difficulties users are struggling to cope with are often  
44 stigmatised, professionals additionally facilitate information flow, fighting social and health  
45 marginalisation or exclusion (see Long et al., 2013). Accordingly, moderators and councillors  
46 play a key role in aiding (peer and professional-user) cooperation as well as preventing abusive  
47 or dangerous behaviours (see Grimmelmann, 2015) perpetrated or suffered by the victims they  
48 support.  
49  
50  
51  
52  
53  
54  
55

56 For the purposes of this paper, we consider ‘responsible AI’ as any AI system which follows  
57 the UKRI framework for responsible innovation. This includes the key principles of Anticipate,  
58

---

59  
60 <sup>1</sup> <https://www.tas.ac.uk/safespacesnlp/>.  
61  
62  
63  
64  
65

1 Reflect, Engage, Act (AREA) and will often involve some sort of stakeholder engagement or  
2 co-design to consider responsibly the environment and context in which the AI system will be  
3 deployed (UKRI, 2022). We focus on Kooth.com<sup>2</sup>, the UK’s leading free online confidential  
4 service (active all year, in the afternoon and evening), which offers counselling and emotional  
5 wellbeing support to young people in the UK. Through its digital platform, users can browse  
6 through self-help materials, seek support or advice on a range of sensitive topics (from bullying  
7 to dealing with suicidal thoughts), share their experience through moderated forums, track their  
8 thoughts and feelings through personal journals, and access synchronous and asynchronous  
9 text-based chats and drop-in sessions with counsellors or emotional wellbeing practitioners.  
10

11  
12  
13 After a brief critical overview of the literature that has looked at moderation and digital  
14 counselling supporting users exhibiting risk behaviours, and a section on our data collection  
15 and analytical approach, this article offers a descriptive account of the working practices at  
16 Kooth Plc (the wider organisation), focusing both on ‘what works’ and on the main challenges  
17 encountered. Departing from these findings, in the conclusions we discuss the possibilities that  
18 a responsible AI can offer to overcome these challenges, without losing track of the positives  
19 in place. We finally signpost where some of the latest trends in responsible AI today might  
20 offer pathways for researchers and practitioners to overcome these challenges.  
21  
22  
23  
24

### 25 26 **Moderation and digital counselling to support users exhibiting risk behaviours**

27  
28  
29 As mentioned before, research on the use of digital platforms for mental health and wellbeing  
30 purposes is not new, as both practitioners, researchers, and even policy makers have recognised  
31 the potential of these digital tools to support users exhibiting risk behaviours (e.g., self-harming  
32 and eating-related disorders), across all ages (Moessner and Bauer, 2012; de la Harpe et al.,  
33 2019; Zhou et al., 2021). For instance, with specific reference to young people – as those  
34 targeted by the services at the centre of our analysis –, recognising that cyberspace has become  
35 a space in which we express ourselves, shape our self-identity, build meaningful relationships  
36 and learn (and hence is a space intrinsically linked to our mental health), the Royal Society for  
37 Public Health (2017) called for action to promote the positive aspects of social media for young  
38 people. This includes access to other people’s health experiences and expert health information,  
39 emotional support and community building, self-expression and self-identity, whilst mitigating  
40 the potential negatives (such as anxiety and depression, sleep, body image, cyberbullying, and  
41 ‘fear of missing out’).  
42  
43  
44  
45  
46  
47  
48

49 Users who typically look for support on these digital platforms can be considered as ‘vulnerable  
50 publics’, as they cohere around socially sensitive and affective issues or experiences, which are  
51 often stigmatised (see McCosker and Wilken, 2017). In this context, the role of moderators in  
52 online communities and digital counsellors (frontline service workers) is of the utmost  
53 importance, as they operate in the blurred lines between caring, or health service work, and  
54 ‘feeling management’ (Hochschild, 2003), in a peculiar type of affective, emotional and  
55 immaterial labour (e.g., McCosker and Darcy, 2013; McCosker, 2018). Overall, these actors  
56  
57  
58  
59

---

60 <sup>2</sup> <https://www.koothplc.com/>.

1 act as brokers facilitating information flow, avoiding marginalisation, social and health  
2 exclusion, and stigma (Long et al., 2013), as they sustain online communities, help framing  
3 and reframing difficult lived experiences, and create and maintain a bridge between the user-  
4 base and professionals (McCosker, 2018). As such, they have a very complex role, as they need  
5 to maintain authority and be perceived as authentic, whilst creating and maintaining trust  
6 (McCosker, 2018).  
7

8  
9 Moderation can be defined as the ‘governance mechanisms that structure participation in a  
10 community to facilitate cooperation and prevent abuse’ (Grimmelmann, 2015: 6). Additionally,  
11 how a group is moderated can influence members’ participation, including their creation and  
12 maintenance of commitment to the community (Ley, 2007). Moderation can take different  
13 forms (West, 2018; Seering, 2020). For the scope of this study, the difference between  
14 automatic or manual moderation matters. Over the years, automated ways to moderate social  
15 media (ranging from classification and filtering approaches, used for instance to identify hate  
16 speech, to more complex digital tools supporting moderation by considering the context of  
17 longer conversations, see e.g. Kurrek et al., 2020; Price et al., 2020; Röttger et al., 2021) have  
18 been developed by social media companies, mostly through AI tools, with the intent of  
19 removing potentially harmful content more effectively and quickly (see, for instance, Gorwa  
20 et al., 2020; Lim et al., 2020). These algorithmic moderation systems have been mostly  
21 analysed and assessed with reference to mainstream social media platforms, fuelled by growing  
22 public expectations for increased platforms’ responsibility. Overall, these systems are often  
23 criticised as being opaque and scarcely effective in complex sociotechnical contexts (Gorwa  
24 et al., 2020), as it can be very difficult for automated tools to make contextual decisions on  
25 complex and multifactorial concepts (Li and Williams, 2018). Also, manual moderation does  
26 come with challenges. For instance, moderation is often carried out by freelancers in poor  
27 working conditions and exposed to extreme amounts of toxic content (Gillespie, 2018)<sup>3</sup>. As  
28 noted, these considerations stemming from moderation research mostly come from analyses of  
29 mainstream platforms. Therefore, a research gap has been identified in considering the realities  
30 and needs of *ad hoc*, more specialised, platforms, such as those focusing on providing services  
31 for mental health and wellbeing.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

43 In digital platforms focusing on mental health and wellbeing, moderation is often sided by  
44 different forms of counselling (e.g., moderators prompting at-risk users to access counselling  
45 services, or counsellors being active in online moderation). Digital counselling, despite its  
46 increasing popularity, is a service still considered complex and controversial (Hendry et al.,  
47 2017; Saha, 2020; Stoll et al., 2020; Perry et al., 2021; Barker and Barker, 2022; Khan et al.,  
48 2022). In digitally mediate service encounters, counsellors deal with relatively new challenges,  
49 mostly linked to the increased accessibility and participation to the services offered, but also  
50 linked to the type of interactions they come across (which entail, for instance, reduced  
51 emotional proximity and the absence of non-verbal cues, see Bambling et al., 2008), their  
52 broader administrative tasks (Tummers et al., 2015; Breit et al., 2021), and risks of vicarious  
53 traumatisation (Furlonger and Taylor, 2013).  
54  
55  
56  
57  
58

---

59  
60 <sup>3</sup> It is noted that this is not the case for the organisation targeted in this study.  
61  
62  
63  
64  
65

1 Moderators and digital counsellors need both intellectual and social capital, as they require  
2 both specialised subject knowledge and the ability to navigate online support in effective ways  
3 (Mowlabocus et al., 2015, as discussed also in McCusker, 2018). In this context, intellectual  
4 capital (subject knowledge) mainly refers to expertise in mental health. As a multidimensional  
5 concept, social capital can be defined as the connections among individuals and social  
6 networks, and the norms of reciprocity and trustworthiness that arise from them (see Putnam,  
7 2000:19). Being the glue holding together social collectives (Sum et al., 2008), social capital  
8 can facilitate the resolution of cooperative action problems (Coleman, 1988; Putnam et al.,  
9 1994). At the core of this concept, is the idea that there are abilities and values rooted in social  
10 networks and relationships, and that these can be achieved through investment in social  
11 relationships; unlike the other forms of capital, no individual ‘owns’ these abilities and values,  
12 as they are only created through interactions across social networks (as summarised in Sum et  
13 al., 2008).

14 While recognising the pivotal role of moderators and digital counsellors, in framing mental  
15 health and recovery practices, it is important to recognise how their ability to act in certain  
16 ways is, in turn, framed by social media affordances, as they create possibilities that both enable  
17 and constrain action (Gibson, 1977; Hutchby, 2001; Bloomfield et al., 2010). Indeed, digital  
18 platforms, including those focusing on providing services for mental health and wellbeing, are  
19 best understood as sociotechnical assemblages and complex institutions (Gillespie, 2017), and  
20 can be conceptualised as composite human (users and, depending on the platform, moderators)  
21 and algorithm-driven non-human (automated tools and filters) entities embedded in their users’  
22 general communicative practices (in line with Prochazka, 2019). As such, to fully understand  
23 the role of moderation and digital counselling, but also their challenges and possibilities, we  
24 cannot avoid considering the specific features of the platforms used (for instance, whether the  
25 communication is asynchronous or synchronous, whether it is organised according to threaded  
26 topics or time-based sequences, or the level of anonymity possible), as these aspects can  
27 directly affect communicative patterns and influence community cohesion, with implications  
28 in terms of the self-disclosure of members and their exchanges of social support (as discussed  
29 in Li et al., 2021). For instance, it has been suggested that, in synchronous communications  
30 (e.g., live chats), members of the community can communicate faster and, thus, form tighter  
31 connections; also, timely feedback seems to play a core role in fostering attachment between  
32 members, probably as speed works as a cue signalling support (Li et al., 2021). As such, we  
33 cannot ignore the importance of social media affordances in the context of social capital and,  
34 specifically, commitment, as digital spaces both enable and constrain certain behaviours,  
35 interactions, and even forms of thought (Ley, 2007).

36 In what follows, we present our study, which furthers research on moderation and digital  
37 counselling to support users exhibiting risk behaviours by looking at the specific context of a  
38 specialised service, offering digital counselling and emotional wellbeing support to young  
39 people in the UK. In doing so, we explored the practices and perceptions of key actors  
40 (moderators, counsellors and individuals in key managerial roles), particularly in relation to  
41 the main challenges moderators face when performing their roles, with a specific attention to  
42

1 the identification of potentially risk behaviours, in the conversations where they provide  
2 support.  
3  
4

## 5 **Methodology** 6

7  
8 We interviewed a total of 6 Kooth.com’s Emotional Wellbeing Practitioners (tiered **across**  
9 trainees **to** more experienced individuals doing moderation and other agreed emotional  
10 wellbeing support with users), 3 Counsellors (who have a clinical and therapy accreditation by  
11 a professional body, and also perform moderation especially in high-risk situations), and 3  
12 Subject Matter Experts (with responsibilities towards the community and its moderation or  
13 focusing on research and operations). Respondents have been identified in the article as PTSs  
14 1-12, see Table 1.  
15  
16  
17  
18

19 Interviews took place through 4 individual interviews (**with PTS6, PTS7, PTS8 and PTS9,**  
20 **respectively**) and 3 focus groups (**firstly, with PTS1, PTS2 and PTS3, secondly with PTS4 and**  
21 **PTS5, and lastly with PTS10, PTS11 and PTS12**) carried out in October and November 2021,  
22 and in March 2022. **Convenience sampling was used during the recruitment. Accordingly,**  
23 **interviews and focus groups were scheduled to suit participants’ availability. Since not all**  
24 **Kooth professionals could take part in the research on the same date and at the same time, both**  
25 **interviews and focus groups were conducted. Furthermore, to ensure that the data collection**  
26 **process was as efficient, smooth and convenient as possible, participants were recruited with**  
27 **the support of the organisation's Research and Evaluation Lead, who advertised the opportunity**  
28 **to engage voluntarily. Moderators were facilitated to attend during their work.**  
29  
30  
31  
32  
33  
34

35 The in-depth interviews were carried out online through the platform Teams, and video-  
36 recorded to keep track also of non-verbal cues. The audio (for a total of 5.30 hours) was then  
37 transcribed and anonymised, in line with the procedure approved by the Ethics Committee of  
38 the University of Southampton (ethical approval ref ERGO/FEPS/66387). **The interviews and**  
39 **focus group discussions were semi-structured, thus followed a pre-defined guide, on the basis**  
40 **of the project’s research questions. Slides with key queries were shared with the participants to**  
41 **facilitate their flow, and to remind participants what they were asked<sup>4</sup>.**  
42  
43  
44  
45  
46  
47  
48

---

49 <sup>4</sup> First, each participant was asked to introduce themselves. Specifically, they had to comment about  
50 their role within Kooth, average workload and challenges frequently faced in their jobs, and the nature  
51 of their roles. This first part of the plan was meant to provide a background to the professional role of  
52 the participants. Then, a series of questions focussed on risk behaviours they had to identify, moderation  
53 practices and forms of interventions, temporal sequencing of actions, professionals involved, reporting  
54 or record-keeping, temporary urgency and strategies typically used to respond to it, as well as indicators  
55 Kooth councillors and moderators tended to look for in young users they supported online. Finally,  
56 specific examples were asked, which served to provide supporting evidence and clarity to the points  
57 made by the participants.  
58  
59  
60  
61  
62  
63  
64  
65

Table 1 – Interviewees

Interviewee	Role at Kooth
PTS1	Emotional Wellbeing Practitioner
PTS2	Emotional Wellbeing Practitioner
PTS3	Subject Matter Experts
PTS4	Emotional Wellbeing Practitioner
PTS5	Emotional Wellbeing Practitioner
PTS6	Counsellor
PTS7	Emotional Wellbeing Practitioner
PTS8	Counsellor
PTS9	Emotional Wellbeing Practitioner.
PTS10	Subject Matter Experts
PTS11	Counsellor
PTS12	Subject Matter Experts

All the transcribed material was manually coded (directed content analysis – see Hsieh and Shannon, 2005) according to the coding scheme summarised in Table 2 (with codes identified *a priori*, in light of our research aim, and subcodes partially adjusted throughout the analysis), and then organised in the following main themes: roles and responsibilities; risks; what works; and current challenges.

Table 2 – Coding framework

CODES	SUBCODES
The actor	Self-definitions; Previous/parallel experience; Background; Tasks; Challenges; Training received; Support received
The work	Specialisation; Shifts; Team; What works; What can be improved and how; Numeric indications (users to deal with/shift; submissions/quarter, etc.); Stages of the work; Rating system
The platform	What works; What can be improved; Potential issues with (semi)automatization
Risk behaviours	Physical health; Self-harm; Sexual health; Mental health; Other
Other	Emerging issues/proxy indicators of problems; COVID-19; Links to criminal behaviour/gangs

Two researchers (Author 1 and 2) undertook the coding ensuring inter-coder consistency (Sanders and Cuneo, 2010; O’Connor and Joffe, 2020), while the project PI (Author 3) was involved in the writing-up of the results and contributed to the refinement of this contribution.



1 Any difference in interpretation among the researchers was addressed through discussions, and  
2 clarifications – if needed – were sought by engaging with the Research and Evaluation Lead at  
3 Kooth Plc. For practicality, the coding procedure was conducted manually, using a Word  
4 document on which every relevant portion of transcription was highlighted, according to a  
5 previously agreed colour-coding scheme. Whilst colours were used to indicate codes,  
6 comments were employed to signal subcodes. This unorthodox strategy was chosen after a  
7 discussion with other members of the multidisciplinary research team (to be involved in other  
8 stages of the project) and Kooth, as it allowed researchers from different backgrounds, at times  
9 non-familiar with qualitative research and coding strategies, to access and monitor the  
10 annotated dataset without having to access specialized software.  
11  
12  
13  
14

15 While extensive reflections on the benefits and the challenges of working in multidisciplinary  
16 research teams, and in having representatives of the organization object of the analysis as part  
17 of the broader research team in the underlying project (*SafeSpacesNLP* – see in the  
18 Introduction) would exceed the scope of this contribution, it is important to briefly underline  
19 how these aspects had a direct impact on our research design. First of all, it is important to note  
20 that, in order to avoid biases, as regards the study presented in this contribution, Kooth’s  
21 representatives were involved only as gatekeeper to facilitate access to potential respondents,  
22 and in helping the researchers to clarify some aspects regarding organizational aspects at  
23 Kooth; no direct input was given on the data collection or analysis process. However, because  
24 of Kooth involvement in other stages of the project, the Research and Evaluation Lead at Kooth  
25 was able to provide constructive feedback on the research design of the broader project, and he  
26 was involved in the ethical oversight of the project and compliance with our data sharing policy  
27 throughout.  
28  
29  
30  
31  
32  
33

## 34 **Results**

### 35 *Roles and responsibilities*

36  
37  
38 From the interviews, the complexity, sensitivity and fluidity of the roles of Emotional  
39 Wellbeing Practitioner, Counsellor and Subject Matter Expert clearly emerged. First of all,  
40 their tasks are performed in a multi-platform and multi-layered environment (referred to as ‘the  
41 platform’ in the article), with relevant information coming from (synchronous and  
42 asynchronous) chats, direct messages, a users’ forum comprising a discussion board and  
43 articles (whose posts and comments are moderated), users’ personal journals and written goals,  
44 and a service inbox. Additionally, there is a dedicated instant messaging channel for staff to  
45 exchange information and get help and clinical support from colleagues and senior shift leads,  
46 in what has been described consistently as a ‘*nurturing environment*’ (PTS2). In order to keep  
47 track of everyone’s work, a dedicated spreadsheet is used (‘*so that we’re not stepping on each  
48 other’s toes and not all clustering is in one area*’, PTS2) to record the team and service tasks.  
49 In this way, moderators and counsellors can focus on certain textual information or on specific  
50 users (some users, for instance, have a named worker allocated to them), in line with their  
51 seniority, when there are at-risk situations.  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 Some tasks have been especially difficult to quantify for interviewees, as they can change every  
2 day. These included the number of submissions they had to address, which depending on level  
3 of risk and complexity could vary from 10 to 100 per shift. Similarly, while caseloads are set  
4 for 1:1 chats, other types of moderation can be more fluid, for instance when moderating live  
5 a discussion board (a type of service-led discussion forum scheduled at a specific time). The  
6 fluidity of the work depends on the specific features of the digital services available which can  
7 offer a combination of synchronous and asynchronous communication. For example, while  
8 users can submit their journals or send comments in forum discussions 24/7, the synchronous  
9 chat with practitioners is open at specific times only. Depending on the role and the seniority,  
10 alongside counselling and moderation, time is allocated to specialised training or to  
11 administrative and managerial tasks.  
12  
13  
14  
15

16 Interactions with users are strictly regulated clinical governance processes, both as regards  
17 access to certain services (particularly live chats with Counsellors and Emotional Wellbeing  
18 Practitioners), and as regards content. For instance, there is a shared resource document with  
19 *'what we're allowed to send to young people, so it's a big list of, like, websites and NHS*  
20 *resources, different organisations, that we can kind of send links to'* (PTS7); specific goals  
21 (e.g., contacting their GP) are to be set or reviewed in each chat; or end-of-chat messages are  
22 to be sent as an encouragement, a recap of the session, or reminder of any goals set or helpful  
23 strategies discussed.  
24  
25  
26  
27  
28

29 A significant part of the moderators' task concerns enabling forms of peer support (*'A lot of*  
30 *the time they'll talk about things that aren't risky though, and that's about maintaining their*  
31 *wellbeing another way. So they might talk about their favourite TV show or something, and in*  
32 *that case we are probably not interacting with them, just kind of publishing that post and*  
33 *facilitating them getting that peer support'* (PTS3)) and making sure that the content shared  
34 with other users is appropriate (for instance, moderators might have to edit posts and comments  
35 keeping them *'as close to how the young persons express themselves as possible but de-*  
36 *escalating risk, [...] maybe deleting some bits and then publishing it'* while interacting with the  
37 user privately (PTS1)). Content moderation is guided by users' age: there are different age  
38 ratings and, as explained by the respondents, *'there shouldn't be any interaction between those*  
39 *age ranges to keep people safe obviously, and the kind of things that people are discussing, the*  
40 *life experience is very different'* (PTS9); *'what may be suitable for [some] might not suitable*  
41 *for [younger people], so we've got to double check'* (PTS1).  
42  
43  
44  
45  
46  
47  
48

49 Depending on the level of risk evidenced in the conversations, *'the engagement levels change'*  
50 (PTS1) users might be referred to a specialist (*'As moderators we don't tend to do the deeper,*  
51 *the therapeutic side of things [...] We tend to focus more on trying to get them into the team,*  
52 *you know, like to the counsellors'* (PTS6)). Counselling takes place digitally and is text-based  
53 (even if, in some parts of the country, there is the possibility to access face-to-face counselling).  
54  
55  
56

57 The workforce employed is pluralistic, with expertise in mental health for young people,  
58 coming both from specialised educational backgrounds (many are qualified counsellors, or  
59 counsellors in training) and from diverse types of mental health practical experience (e.g., *'I*  
60  
61  
62  
63  
64  
65

1 have actually been moderating mental health communities online for about 20 years outside of  
2 Kooth' (PTS1); 'I've done a lot of work in schools [...] and as a support worker, kind of building  
3 up to doing my counselling qualification' (PTS8)).  
4

5 The complexity and fluidity of the context, not surprisingly, creates a series of role tensions,  
6 with Emotional Wellbeing Practitioners, Counsellors and Subject Matter Experts alike having  
7 to juggle different needs. A main aspect refers to their own wellbeing (as tasks can be  
8 'emotionally draining' (PTS4)):  
9

10  
11  
12 *'It's not the it's not just dealing with the one risk, it's dealing with multiple bits*  
13 *of risk and it all just kind of layering on and then by the end of the day you just*  
14 *like 'Jesus, that was a lot of heavy different topics on a load of different stuff as*  
15 *well' [...] But then there's some that bring you right back and they hit you right*  
16 *in your stomach and you really feel them and they can impact you afterwards'*  
17 (PTS7).  
18  
19  
20  
21

22 All respondents discussed the importance of self-care and of setting their own boundaries, both  
23 by themselves ('I like lighting a candle at eight o'clock. So, lighting a candle, I'm putting my  
24 music on. And that's how I... get through my shift' (PTS5)), or through peer support ('The peer  
25 support we provide each other as colleagues and [...] having awareness of things like burnout  
26 and vicarious trauma [...] I can probably see the word suicide 100 times a day if I'm not careful,  
27 that kind of thing. So kind of prolonged and chronic and that kind of exposure' (PTS3)).  
28 Additionally, those involved in moderation have regular meetings with their line managers and  
29 can access clinical support and external supervision.  
30  
31  
32  
33  
34

35 A second important aspect refers to the need to manage time effectively through complex  
36 situations, especially after performance indicators (used for benchmarking) were introduced in  
37 recent years. That change, in the words of some respondents, 'added pressure' (PTS5), creating  
38 some 'rush' that can be difficult to manage when 'dealing with emotional [...] wellbeing and  
39 trauma' (PTS4).  
40  
41  
42

### 43 **Risks**

44  
45  
46 Not surprisingly, assessing risk is central in moderators and counsellors' activities ('So it's  
47 almost a constant, every minute that we're working we're assessing for risk in one way or  
48 another from multiple directions' (PTS1)), in what has been described as a 'better safe than  
49 sorry' approach (PTS2). During every shift, a couple of moderators **oversaw** scanning  
50 messages for risks, to escalate those needing more urgent attention in the platform.  
51  
52  
53

54 As collectively reported by the respondents, risks in Kooth.com often relate to eating disorders,  
55 anxiety, depression and other types of mental health issues, gender dysphoria, sexual health,  
56 self-harm behaviours, suicide attempts, bullying, physical and mental abuse, sexual  
57 exploitation, forced marriages, and grooming. Also, victims of crime or young people exploited  
58 in crime (for instance, young people involved in gangs for drug dealing) report their  
59  
60  
61  
62  
63  
64  
65

1 experiences online. These risks can also be multiple, and some of those risks worsened during  
2 COVID-19 restrictions (*'we've just seen numbers go through the roof'* (PTS3) - see Gerrard,  
3 2020, on the surge in demand for mental health charities during the pandemic). However, as  
4 reported by one respondent, in many cases the risks identified are *'the early levels, so kind of*  
5 *early emerging eating difficulties [or] we're seeing situations escalate, so maybe it's a situation*  
6 *with bullying that is becoming physical'* (PTS3).  
7  
8

9  
10 There can be cases where a risk has already escalated, or there is an imminent danger (and, as  
11 such, external services are called: *'So if they disclose that they've self-harmed badly, taken an*  
12 *overdose or a severe risk, we can call ambulances for them'* (PTS3)). And some users could  
13 be more at-risk than others. The service has multiple processes in place to quickly communicate  
14 to staff what level of risk someone is currently assessed at, and any key aspects of their care  
15 plan to be aware of.  
16  
17

### 18 19 **What works** 20

21  
22 Overall, the respondents were very positive regarding the social utility of their work, as  
23 evidenced for instance in the following snippet: *'Even if it takes us a bit of time, we're still far*  
24 *faster than unfortunately people like GPs or CAMHS<sup>5</sup> can be at the moment, so we are still*  
25 *providing, you know, waiting list free, essentially, access to not just support, but to be able to*  
26 *flag up issues and then have them kind of escalated and dealt with professionally'* (PTS3).  
27 Even in less at-risk cases, they can give users *'that little bit of the extra confidence, just to kind*  
28 *of make that step for themselves because they've seen other people's personal experiences of*  
29 *that'* (PTS7).  
30  
31  
32  
33

34  
35 The approach used is considered effective, as it allows a good mixture of both peer and  
36 professional support (*'We've seen the amazing support that these kids give each other on that*  
37 *website, it's just absolutely fantastic'* (PTS2); *'We keep that a nice community and nice place*  
38 *for them all to speak'*(PTS6)). And, reportedly, the feedback received from users is very  
39 positive as well.  
40  
41

42  
43 Personal experiences in the organisation were generally valued positively, especially as regards  
44 team support and the possibility to have some build-in flexibility in their tasks (e.g., *'We are*  
45 *an agile little team [...] It is a very robust system, people are very, very supportive and we work*  
46 *in a really collaborative way'* (PTS3).  
47  
48  
49  
50  
51  
52  
53  
54  
55

---

56 <sup>5</sup> General Practitioners (GP), in the UK are doctors who treat all common medical conditions and refer  
57 patients to hospitals and other medical services; Children and Adolescent Mental Health Service  
58 (CAMHS) in the UK are the public services that assess and treat young people with emotional,  
59 behavioural or mental health difficulties.  
60  
61

## Current challenges

In the previous section, titled ‘what works’, we pointed out approaches and practices that responsible AI systems should maintain, foster and further implement. Building on that discussion, in this section we discuss a number of existing challenges and difficulties identified by our respondents, which are of particular interest in the context of this study as they highlight necessary points of intervention. The challenges reported by the respondents were grouped into the following topics: time use and multitasking; reading in between the lines; and more subtle risks.

### *Time use and multitasking*

Time is a scarce resource, and - because of the volume of work - respondents consistently mentioned (the lack of) time as a major issue in carrying out their tasks at their best, leading to the lack of taking sufficient breaks or allocated time to debrief (*‘it was just too much because I couldn’t really process in between chats’* (PTS8)), especially when there are risk situations (for instance, suicidal intention) (*‘a risky situation can take up a whole shift, especially if there is a immediate risk for the user’s safety’* (PTS4)). Moreover, time issues can lead to less effective support, for instance when users are left waiting too long to access digital counselling, or when errors can be made because of necessary multitasking and difficulty of the task (*‘With a high volume of work, there are going to be errors [...] You know, somebody misses, like I say, like a case note, or they might not have edited...’* (PTS9); *‘I’ve moderated a post that [...] it was not against boundaries, but it did kind of push the boundaries a little bit. [...] I was dealing with multiple things at once when I read it I was like ‘that’s that’s within the boundaries, I’ll publish that’, like the wording wasn’t off, but because it was so short, like you missed their underlying tone, ‘cause sometimes it’s not necessarily what’s written there, but it’s kind of the impact that would have on the community as a whole, it’s like an ecosystem at the end of the day we want a positive one and so’* (PTS7)). Of course, clinical auditing processes are regularly carried out to understand, identify and mitigate human error, including processes that take into consideration the platform and the ecosystem of services that are delivered.

As explicitly discussed by some respondents, a system to help them navigating the mole of information they need to go through would be welcomed, as exemplified in the following snippets:

*‘At the moment a lot of our processes are very manual and us literally just sitting there and reading through kind of really large chunks of text, so I don’t know something that made that those couple of words [e.g., suicide] pop out a bit so they don’t get missed and are easier to pick up’*(PTS3).

*‘Time [...] is a big challenge [...] ‘cause sometimes you’re trying to, you know, keep all the immediate posts, you know, spend your time on them, but then also you’ve got [...] the lower risk posts which are just as important because they want to be heard, they need the peer support [...] They sometimes might be*

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

*[temporarily] left because risk always comes first [...]. So I don't know if there could be a tool to identify, you know, high risk posts' (PTS6).*

### *Reading in between the lines*

Identifying risk is not always straightforward, as users might be less 'direct' in expressing their feelings and emotions, and there are no visual or non-verbal cues to be observed (there are, however, text-based cues). For instance, users could use metaphors (PTS4), or post poetry that needs to be interpreted (PTS2). As such, both moderators and counsellors need to try and slowly build the picture of what is going on, especially since users can sip the amount of information they provide ('*they[users] are definitely in charge*' (PTS4)), and, online, many important risk-relevant cues used in other contexts (such as the body language, the tone of the voice) are missing ('*When I first started doing chats, I used to have a massive headache all the time [...] Swapping over from doing face to face interventions, or even telephone counselling or telephone counselling skills is like massive*' (PTS8)). Consider the following example:

*'A young person could say I have been feeling very low lately... I've got a lot going on in my life, I'm feeling scared, I'm feeling alone. [...] Straight away as moderator I would think: well, why are they feeling scared, why do they feel alone, what's going on in their life?. That could be anything, [...] there could be risk there [...] As moderators we go digging, we want to find out a lot more [...] You can open a can of worms unknowingly' (PTS6).*

### *More subtle risks*

Finally, there are some moderation and counselling challenges that are linked to what could be considered more subtle forms of risk (to the individual user, or to the digital community), and that can complicate the work of employed staff, or hinder the inclusiveness of the service provided. Some have to do with organisational issues (because of the resources available): for instance, despite the extensive operational hours of the services offered at Kooth<sup>6</sup>, it was lamented by one of our respondents that Emotional Wellbeing Practitioners and Counsellors are not available overnight and are available only at limited times over the weekends, which might discourage some potential users from participating, as they might have less time during the rest of the week (PTS8).

More challenges are linked to the content posted, as extra attention is constantly needed to make sure not only that shared content is appropriate, but also appropriate to a specific age group ('*It's acceptable for a 16 year old on our site to ask where they can get free condoms, but if an 11 year old is asking that question we are reacting in a very different way*' (PTS3)). Additionally, attention is constantly needed to ensure that the digital platforms do not foster any type of dis- or mis-information ('*We are very aware we're not medical professionals and*

---

<sup>6</sup> Kooth operates out of hours service weekdays from 12 to 10 pm and the weekend from 6-10 pm 365 days a year.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

*their peers certainly aren't either so with those questions what we tend to do is we'll say: we're really sorry we can't publish this [...] but hey let's have a talk about it, maybe here's an NHS information page we can give you that you can read and come back if you have any questions'* (PTS3)). Also, anonymity and confidentiality **must be** maintained throughout, so that extra effort is needed to check, for instance, that the same username is not used across social media as that could easily reveal the real identity of a user (*'Something that could be useful to do with some automation if possible'*, PTS1). At times, the content posted could be difficult to interpret because slang (from different geographical areas) is used, or users make references to numbers (for instance, of county lines offences) and moderators could not be familiar with those (PTS2).

A final set of risks refers to the need and importance to foster a climate of trust with users, and to build and maintain a relationship (*'you need to build that rapport with them first, before they feel that they can open up'* (PTS5)). This entails showing that staff care about users, but also that they need to be treated with respect (*'You know, we sometimes get asked "are you robots, like are you real people?"'* (PTS3); *'You do get some users who sign up and they just send like erm, ridiculous things. And, or they think that we're like robots. [...] Yeah, or like those and rude things, or.....very inappropriate things. [...] We always message them in a way that shows that we are human, [...] a lot of the time, the users who send something silly, they're just kind of testing the water'* (PTS9)).

## Discussion and conclusions

Digital platforms - as intermediaries bringing together users, service providers, content producers and other stakeholders for a range of social exchanges (Srnicsek, 2017) - are of increasing social importance, to the point that it has been claimed that we now live in a 'platform society' (van Dijck et al. 2018). This society, however, is generally dominated by large scale, monopolistic platforms, and so is most research on digital moderation. In the study presented here, we focused on the contrary on the realities and needs of a specialised platform focusing on providing services for mental health and wellbeing to young people.

As discussed previously, the main goal of the study reported in this article was to further research on moderation and digital counselling to support users exhibiting risk behaviours by looking at the practices and perceptions of key actors involved in moderation, counselling, or in managerial roles. We were particularly interested in the main challenges those involved in substantial moderation practices (in our case, Emotional Wellbeing Practitioners and Counsellors) face when performing their roles, to create a benchmark to ideate and develop ways to improve the system currently in place, exploring ways in which responsible AI solutions can support and better target interventions for young people asking for help or engaging in risk behaviours online.

We have seen that Kooth effectively combines individual counselling with community support and is subject to rigorous ethical standards. If we are to address the thorny issue of safe moderation within digital platforms, it is vital that we learn from well-established services in

1 the field to specifically understand the key challenges and barriers to safe, scalable moderation.  
2 To summarise, the challenges identified by counsellors and moderators we interviewed mainly  
3 revolved around three areas: time management, interpretation of user communication, and  
4 subtle risks. The first category of challenges involved the time constraints imposed by the  
5 limited resources available, despite the number of young people turning to the online platform  
6 for support. Such time challenges were exacerbated by the type of support provided and the  
7 difficulties involved by reducing the time dedicated to the discussion of sensitive topics with  
8 vulnerable users needing assistance. At the same time, these challenges also caused  
9 professionals to limit their time to filter the information provided by the users, and to recover  
10 after these difficult discussions. The second group of challenges regarded how users  
11 communicate on the platform and with the professionals available on it. Since struggles and  
12 needs are often communicated indirectly, calls for support and their urgency **are** not  
13 immediately clear from the user interactions and frequently require mental health workers to  
14 interpret the texts received through the platform and to collect more information about  
15 individuals who authored them, to better understand their communication styles and  
16 preferences. Given the sensitivity of the topics and risks, these tasks require more time and  
17 effort from the moderators and counsellors to process relevant data and accordingly respond to  
18 users. The last set of challenges mentioned during the interviews dealt with the need  
19 professionals have to build and maintain a trustworthy relationship with the young people they  
20 support to be able to help them to the best of their abilities. At the same time, workers also  
21 needed to ensure the suitability of all contents that are publicly shared through the platform,  
22 according to viewers age, vulnerability, and cognitive capacities. Therefore, all the challenges  
23 discussed by counsellors and moderators are intertwined and interrelated.  
24  
25  
26  
27  
28  
29  
30  
31  
32

33 Recent advances in responsible AI methods are now providing insights into possible solutions  
34 to some of these challenges. With regards the challenge of time use and multitasking, modern  
35 text classification algorithms (Minaee et al. 2021) can classify large volumes of online posts to  
36 allow moderators to better triage and filter posts, and to flag time critical posts for example  
37 those that could represent a threat to life such as posts with markers associated with suicidal  
38 behaviour or **immediate** need. The use of AI to organise content can also provide useful  
39 structure **to moderator's** daily workload, reducing the chance of human error due to missed or  
40 forgotten content. The use of AI to classify ambiguous, subtle or contextual age-appropriate  
41 content has led to various grand challenges within the AI research community (Joshi 2017),  
42 but AI approaches that can embed knowledge graphs containing common sense and/or domain-  
43 specific contextual knowledge into AI models may yield results in the medium term. Recent  
44 research projects are exploring a range of novel AI methods to progress text classification  
45 research. For example, SafeSpacesNLP<sup>7</sup> and ProTechThem<sup>8</sup> are exploring AI models that can  
46 move beyond single post classification, such as hate speech classification, and towards an  
47 ability to identify moments of change within conversations around mental health issues. Other  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58

---

59 <sup>7</sup> <https://www.tas.ac.uk/safespacesnlp/>

60 <sup>8</sup> <http://www.protechthem.org/>



1 projects such as the UKRI TAS Trust Node<sup>9</sup> and EU H2020 project WeVerify<sup>10</sup> are exploring  
2 AI models for trustworthiness and misinformation detection. A recent trend is the emergence  
3 of new human-in-the-loop AI approaches (Middleton et al. 2022), which represents an exciting  
4 direction of travel as ultimately online content moderation is a human process and we need AI  
5 algorithms which can be trustworthy and supportive to empower human decision makers to  
6 concentrate more time on the subtle and tricky moderation decisions which AI algorithms find  
7 hard to process, leading to an use of automation that can enhance human activities in  
8 moderation to preserve a safe space and positive digital mental health community.  
9

10  
11  
12 **As such, as responsible AI approaches are becoming increasingly promising in their capacity**  
13 **to support up-scaling of existing manual moderation approaches and better target interventions**  
14 **for vulnerable users in sensitive settings, whilst developing in ways that are ‘fit for purpose’**  
15 **and minimise potential biases.** Sociotechnical approaches bringing together computational  
16 expertise with social sciences and subject matter experts' ability to investigate and interpret  
17 qualitative datasets **is becoming** essential. Looking forward, in our opinion it is the combination  
18 of human experience, with its capacity for insight, connection and empathy, alongside  
19 responsible AI, with its capacity for processing content at scale, that will deliver the step  
20 changes needed to address the significant challenges we have identified for **digitally scaled-up**  
21 moderation in a digital mental health community. If AI is delivered in a responsible way, with  
22 safeguards in place as part of a wider and evidenced governance framework, coupled with  
23 opportunities for stakeholder trust to be built via mechanisms such as codesign, then it can  
24 potentially empower moderators and lead the field of digital community moderation into an  
25 exciting future.  
26  
27  
28  
29  
30  
31

## 32 33 34 35 **Acknowledgements**

36  
37 This work was supported by the Engineering and Physical Sciences Research Council  
38 (EP/V00784X/1) and Economic and Social Research Council (ES/V011278/1). We would like  
39 to thank the Research and Evaluation Lead at Kooth Plc., Dr Santiago De Ossorno Garcia, for  
40 his precious insights, Toni Mees and Dr Hannah Wilson from the clinical team at Kooth to  
41 provide with the time and coordinator for the project, Aaron Sefi and Dr. Lynne Green to help  
42 enabling the project at TAS, and Aynsley Bernard support in recruitment..  
43  
44

## 45 46 **References**

47  
48 Bambling M, King R, Reid W and Wegner K (2008) Online counselling: The experience of  
49 counsellors providing synchronous single-session counselling to young people. *Counselling*  
50 *and Psychotherapy Research* 8, doi: 10.1080/14733140802055011.  
51

52  
53 Barker GG and Barker EE (2022) Online therapy: lessons learned from the COVID-19 health  
54 crisis. *British Journal of Guidance & Counselling* 50(1):66-81.  
55  
56  
57

---

58  
59 <sup>9</sup> <https://trust.tas.ac.uk/>

60 <sup>10</sup> <https://weverify.eu/>  
61  
62  
63  
64  
65

1 Bloomfield BP, Latham V and Vurdubakis T (2010) Bodies, technologies and action  
2 possibilities: When is an affordance? *Sociology* 44(3):415-433.

3  
4 Breit E, Egeland C, Løberg IB and Røhnebæk MT (2021) Digital coping: How frontline  
5 workers cope with digital service encounters. *Soc Policy Adm.* 55:833-847.

6  
7 Coleman J (1988) Social Capital in the Creation of Human Capital. *American Journal of*  
8 *Sociology* 4:95-121.

9  
10 de la Harpe R, Settley C and Cilliers R (2019) Online counselling services for Youth@risk.  
11 *CONF-IRM 2019 Proceedings*, 32.

12  
13 Furlonger B and Taylor W (2013) Supervision and the Management of Vicarious  
14 Traumatization Among Australian Telephone and Online Counsellors. *Australian Journal of*  
15 *Guidance and Counselling* 23(1):82-94.

16  
17  
18  
19 Gerrard Y (2020) The COVID-19 Mental Health Content Moderation Conundrum. *Social*  
20 *Media + Society*. doi:10.1177/2056305120948186

21  
22  
23 Ghallab M (2019) Responsible AI: requirements and challenges. *AI Perspect* 1(3).

24  
25 Gillespie T (2018) *Custodians of the Internet: Platforms, Content Moderation, and the Hidden*  
26 *Decisions that Shape Social Media*. New Haven, CT: Yale University Press.

27  
28  
29 Gibson J (1977) The theory of affordances. In: Shaw R and Bransford J (eds.) *Perceiving,*  
30 *Acting, and Knowing: Toward an Ecological Psychology*. Hillsdale, NJ: Erlbaum, 67-82.

31  
32  
33 Gorwa R, Binns R and Katzenbach C (2020) Algorithmic content moderation: Technical and  
34 political challenges in the automation of platform governance. *Big Data & Society*,  
35 <https://doi.org/10.1177/2053951719897945>.

36  
37  
38 Grimmelmann J (2015) The virtues of moderation. *Yale Journal of Law & Technology* 17:42.

39  
40 Hansen MC and Aranda MP (2012) Sociocultural influences on mental health service use by  
41 Latino older adults for emotional distress: Exploring the mediating and moderating role of  
42 informal social support. *Social Science & Medicine*, 75(12):2134-2142.

43  
44  
45 Hendry NA, Robards B and Stanford S (2017) Beyond social media panics for ‘at risk’ youth  
46 in mental health practice. In: Stanford S, Sharland E and Heller NR (eds.) *Beyond the Risk*  
47 *Paradigm in Mental Health Policy and Practice*. Basingstoke: Palgrave Macmillan, pp.135-  
48 154.

49  
50  
51 **Hsieh HF and Shannon SE (2005) Three approaches to qualitative content analysis.**  
52 ***Qualitative health research*, 15(9):1277-1288.**

53  
54  
55 Hochschild AR (2003) *The Commercialization of Intimate Life*. Berkeley, CA: University of  
56 California Press.

57  
58  
59 Hutchby I (2001) Technologies, texts and affordances. *Sociology* 35(2):441-456.

1 Joshi, A. Bhattacharyya, P. Carman, M.J. (2017) Automatic Sarcasm Detection: A Survey.  
2 ACM Comput. Surv. 50, 5, Article 73 (September 2018), 22 pages.  
3 DOI:<https://doi.org/10.1145/3124420>

4  
5 Khan S, Shapka JD and Domene JF (2022) Counsellors' experiences of online  
6 therapy. *British Journal of Guidance & Counselling* 50(1):43-65.

7  
8 Kivitz J (2013) E-Health and renewed sociological approaches to health and illness. In:  
9 Orton-Johnson K and Prior N (eds.) *Digital Sociology: Critical*  
10 *Perspectives*. Basingstoke: Palgrave Macmillan, pp, 213-226.

11  
12 Kurrek J, Saleem HM and Ruths D (2020) Towards a Comprehensive Taxonomy and Large-  
13 Scale Annotated Corpus for Online Slur Usage, Workshop on Online Abuse and Harms,  
14 EMNLP.

15  
16  
17  
18 Lederman R, Fan H, Smith S and Chang S (2014) Who can you trust? Credibility assessment  
19 in online health forums. *Health Policy and Technology*, 3(1):13-25.

20  
21  
22 Ley BL (2007) Vive Les Roses!: The Architecture of Commitment in an Online Pregnancy  
23 and Mothering Group. *Journal of Computer-Mediated Communication* 12(4):1388-1408.

24  
25  
26 Li H, Kraut RE and Zhu H (2021) Technical Features of Asynchronous and Synchronous  
27 Community Platforms and their Effects on Community Cohesion: A Comparative Study of  
28 Forum-based and Chat-based Online Mental Health Communities, *Journal of Computer-*  
29 *Mediated Communication* 26(6):403-421.

30  
31  
32 Li S and Williams J (2018) Despite What Zuckerberg's Testimony May Imply, AI Cannot  
33 Save Us. *Electronic Frontier Foundation Deeplinks Blog*. Available  
34 at: [https://www.eff.org/deeplinks/2018/04/despite-what-zuckerbergs-testimony-may-imply-](https://www.eff.org/deeplinks/2018/04/despite-what-zuckerbergs-testimony-may-imply-ai-cannot-save-us)  
35 [ai-cannot-save-us](https://www.eff.org/deeplinks/2018/04/despite-what-zuckerbergs-testimony-may-imply-ai-cannot-save-us).

36  
37  
38 Lim Y, Lim CM, Gan KH and Samsudin N (2020) Text Sentiment Analysis on Twitter to  
39 Identify Positive or Negative Context in Addressing Inept Regulations on Social Media  
40 Platforms. *2020 IEEE 10th Symposium on Computer Applications & Industrial Electronics*  
41 *(ISCAIE)*, 96-101.

42  
43  
44 Long JC, Cunningham FC, Braithwaite J (2013) Bridges, brokers and boundary spanners in  
45 collaborative networks: a systematic review. *BMC Health Services Research* 13:158.

46  
47  
48 McCosker A and Darcy R (2013) Living with cancer: affective labour, self-expression and  
49 the utility of blogs. *Information, Communication & Society* 16(8):1266-1285.

50  
51  
52 McCosker A (2018) Engaging mental health online: Insights from beyondblue's forum  
53 influencers. *New Media & Society* 20(12):4748-4764.

54  
55  
56 McCosker A and Wilken R (2017) Mapping mental health intermediaries: vulnerable publics  
57 and platformed support. Paper presented at AoIR 2017: the 18th annual conference of the  
58 Association of Internet Researchers, Tartu, 18-21 October. Available at: <http://spir.aoir.org>.

1 Middleton, S.E. Letouzé, E. Hossaini, A. Chapman, A. (2022) Trust, regulation, and human-  
2 in-the-loop AI: within the European region, *Communications of the ACM (CACM)*, 65, 4  
3 (April 2022), 64–68. DOI:<https://doi.org/10.1145/3511597>

4 Minaee, S. Kalchbrenner, N. Cambria, E. Nikzad, N. Chenaghlu, M. Gao, J. 2021. Deep  
5 Learning--based Text Classification: A Comprehensive Review. *ACM Comput. Surv.* 54, 3,  
6 Article 62 (April 2022), 40 pages. DOI:<https://doi.org/10.1145/3439726>

7 Moessner M and Bauer S (2012) Online counselling for eating disorders: Reaching an  
8 underserved population? *Journal of Mental Health* 21(4).

9 Moorhead SA, Hazlett DE, Harrison L (2013) A new dimension of health care: systematic  
10 review of the uses, benefits, and limitations of social media for health communication.  
11 *Journal of Medical Internet Research* 15(4):e85.

12 Mowlabocus S, Harbottle J and Tooke B (2015) ‘Because even the placement of a comma  
13 might be important’: expertise, filtered embodiment and social capital in online sexual health  
14 promotion. *Convergence* 21(3):375-387.

15 National Institute of Mental Health (2017) Technology and the future of mental health  
16 treatment, February. Available at: [https://www.nimh.nih.gov/health/topics/technology-and-  
17 the-future-of-mental-health-treatment/index.shtml](https://www.nimh.nih.gov/health/topics/technology-and-the-future-of-mental-health-treatment/index.shtml).

18 O’Connor C and Joffe H (2020) Intercoder reliability in qualitative research: debates and  
19 practical guidelines. *International journal of qualitative methods*, 19:1609406919899220.

20 Perry A, Pyle D, Lamont-Mills A, et al. (2021) Suicidal behaviours and moderator support in  
21 online health communities: a scoping review. *BMJ Open* 11:e047905.

22 Price I, Gifford-Moore J, Flemming J, Musker S, Roichman M, Sylvain G, Thain N, Dixon L  
23 and Sorensen J (2020) Six Attributes of Unhealthy Conversations, Workshop on Online  
24 Abuse and Harms, EMNLP.

25 Putnam R, Leonardi R and Nanetti R (1994) *Making Democracy Work: Civic Traditions in  
26 Modern Italy*. Princeton: Princeton University Press.

27 Putnam RD (2000) *Bowling alone: The collapse and revival of American community*. Simon  
28 and Schuster.

29 Röttger P, Vidgen B, Nguyen D, Waseem Z, Margetts H and Pierrehumbert J (2021)  
30 HateCheck: Functional Tests for Hate Speech Detection Models, ACL.

31 Royal Society for Public Health (2017) #StatusOfMind: social media and young people’s  
32 mental health and wellbeing. Available at: [https://www.rsph.org.uk/our-  
33 work/campaigns/status-of-mind.html](https://www.rsph.org.uk/our-work/campaigns/status-of-mind.html).

34 Saha K, Ernala SK, Dutta S, Sharma E and De Choudhury M (2020) Understanding  
35 Moderation in Online Mental Health Communities. In: Meiselwitz G (eds.) *Social Computing  
36 and Social Media. Participation, User Experience, Consumer Experience, and Applications*

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

of *Social Computing*. HCII 2020. Lecture Notes in Computer Science, vol 12195. Cham: Springer.

Sanders CB and Cuneo CJ (2010) Social reliability in qualitative team research. *Sociology*, 44(2):325-343.

Seering J (2020) Reconsidering Self-Moderation: The Role of Research in Supporting Community-Based Models for Online Content Moderation. *Proceedings of the ACM on Human-Computer Interaction* 4:107-128.

Sokol R and Fisher E (2016) Peer support for the hardly reached: a systematic review. *American Journal of Public Health* 106(7):e1-e8.

Srnicek N (2017) *Platform Capitalism*. Cambridge: Polity.

Stoll J, Müller JA and Trachsel M (2020) Ethical issues in online psychotherapy: A narrative review. *Frontiers in Psychiatry* 10, 993.

Sum S, Mathews MR, Pourghasem M and Hughes I (2008) Internet Technology and Social Capital: How the Internet Affects Seniors' Social Capital and Wellbeing. *Journal of Computer-Mediated Communication* 14(1):202-220.

Tanis M (2008) Health-related on-line forums: what's the big attraction? *Journal of Health Communication* 13(7):698-714.

Tucker IM and Goodings L (2017) Digital atmospheres: affective practices of care in Elefriends. *Sociology of Health & Illness* 39(4):629-642.

Tummers L, Bekkers V, Vink E and Musheno M (2015) Coping during public service delivery: A conceptualization and systematic review of the literature. *Journal of Public Administration Research and Theory* 25(4):1099-1126.

UKRI (2022) Framework for responsible innovation. Available at: <https://www.ukri.org/about-us/epsrc/our-policies-and-standards/framework-for-responsible-innovation/>.

Van Dijck J, Poell T and de Waal M (2018) *The Platform Society: Public Values for a Connective World*. Oxford: Oxford University Press.

West SM (2018) Censored, suspended, shadow banned: User interpretations of content moderation on social media platforms. *New Media Soc.*, 20.

Zhou X, Bambling M and Edirippulige S (2021) A mixed-method systematic review of text-based telehealth interventions in eating disorder management. *Journal of Health Research* (online first).

## REBUTTAL LETTER

Dear Editor, dear Reviewers,

We are submitting a revised (and proofread by a native speaker) version of the article entitled *Identifying key challenges and needs in digital mental health moderation practices supporting users exhibiting risk behaviours to develop responsible AI tools: the case study of Kooth* we had previously sent for consideration in SN Social Sciences.

We have addressed the comments by **Reviewer #1** as follows (in red):

- Were descriptive sociodemographics not collected for the sample? It's currently impossible to contextualise the findings demographically beyond job role. **To preserve the anonymity of our respondents – a necessary step to comply with the ethical approval we received for this research study – it is not possible to add more detailed demographic information in the paper. Please note that the number of Emotional Wellbeing Practitioners, Counsellors, and Subject Matter Experts in the organization we studied (and whose name was made public) is limited, so unfortunately any additional information would hinder their anonymity.**
- Were software used to undertake the analysis? It would be standard practice to outline the process of undertaking the analysis. **No software was used to undertake the analysis. We added more details about the steps we undertook in the methodology section (after Table 2).**
- Did the interviews follow a topic guide? Unclear the level of rigour involved in these interviews (were they semi-structured etc). **Yes, the interviews were semi-structured, and thus followed a topic guide. Slides listing the key questions were shared with the participants during the interviews or focus groups to facilitate their flow and remind participants of the questions they were asked. In the revised manuscripts, more details were added before Table 1 (please note also the footnote).**
- How many people undertook the coding? Further, if multiple people were involved how were any differences in interpretation resolved? If only one person coded the data how were the wider study team involved in the analytical process? **Two researchers undertook the coding ensuring inter-coder consistency through extensive discussion on the interpretation of data. The project PI was involved in the writing-up of the results and contributed to the refinement of this contribution, but was not involved directly in the qualitative analysis. These additional specifications regarding the analysis process have been added to the methodology section (see the paragraph after Table 2).**
- What sampling strategy was used during recruitment? **Convenience sampling was used during the recruitment. These additional specifications regarding the analysis process have been added to the methodology section (see the second paragraph).**
- How was data saturation navigated? **We analysed the full dataset (all the interviews and the focus group discussions, in their entirety) as now explicitly specified after Table 1, even when the coding process was revealing of only redundant themes according to both the researchers involved in the analysis.**

- What specific data analysis approach was undertaken and was a citable process followed? It was also unclear how 'domains' may fit into this approach. **Directed content analysis (Hsieh and Shannon, 2005) was conducted, as now explicitly specified in the revised manuscript. In the revised version, we substituted the word 'domains' with 'topics' to avoid potential confusion.**
- Reflexivity may also be especially important, given Kooth staff members such as the Research and Evaluation Lead were on the research team. **We added a few notes on this point in the revised manuscript (before the Result section).**

We have addressed the comments by **Reviewer #2** as follows (in red):

- It is not clear why you used both focus groups and interviews and what data came from what method. **For convenience, interviews and focus groups were scheduled according to participants' availability. Since not all of them could be available on the same date and at the same time, both interviews and focus groups were conducted. These additional specifications regarding the analysis process have been added to the methodology section (see the second paragraph).**
- I would appreciate more details regarding the roles of your participants - explain the specific role in more detail. **The level of detail provided has been negotiated with the organization we studied (because of the sensitive role of our research participants, the organization was not comfortable with us sharing more detailed information, as that might negatively affect their work). The precise wording used to describe the specific roles has been agreed with them. Please also note, as already stressed in the responses to the other Reviewer, that adding more detailed information on specific respondents would have hindered their anonymity.**
- Currently you don't define risk until page 10. I think you need to do this in your introduction, so it is really clear from the start what you are defining as risk. **Several examples of risk behaviours have been added to the fifth paragraph of the introduction. Additionally, risk behaviours are now discussed in more detail in the first paragraph of the section titled 'Moderation and digital counselling to support users exhibiting risk behaviours'.**
- I think you need to explain the role of a moderator in your introduction. Currently on page 4 (paragraph 2) is where you talk about moderators but it is hard to follow and I believe it could be presented in a more coherent way. **The role of moderators has been briefly summarised in the sixth paragraph of the introduction to anticipate this key concept to the reader. To increase the clarity and improve the flow of the paper, this paragraph signposts to the following section, where the moderation role is now discussed more extensively.**
- Ensure you are explaining all your concepts - I am specifically referring to "expert-client relationship" and "public - professional relationship". **These concepts have now been explained in the Introduction.**

We would like to thank both Reviewers for their constructive feedback.

The Authors.