

# Investigating radicalisation indicators in online extremist communities

**Christine de Kock**

University of Melbourne  
christine.dekock@unimelb.edu.au

**Eduard Hovy**

University of Melbourne  
eduard.hovy@unimelb.edu.au

## Abstract

We identify and analyse three sociolinguistic indicators of radicalisation within online extremist forums: hostility, longevity and social connectivity. We develop models to predict the maximum degree of each indicator measured over an individual’s lifetime, based on a minimal number of initial interactions. Drawing on data from two diverse extremist communities, our results demonstrate that NLP methods are effective at prioritising at-risk users. This work offers practical insights for intervention strategies and policy development, and highlights an important but under-studied research direction.

## 1 Introduction

Online extremism is a pressing problem with a proven relation to not only indirect societal harm (Blake et al., 2021) but also to concrete offline dangers in the form of terrorist activities (Gill et al., 2017; Baele et al., 2023). Though disconcerting, the growth of publicly available online content that espouses extremist views presents an opportunity to use computational methods for detecting, channelling, and combating extremist behaviour.

Despite the significance of language to this issue, there has been limited NLP research on extremism and radicalisation. Existing work has focused on behaviours related to specific communities. For instance, de Gibert et al. (2018) introduced a dataset of hate speech on a white supremacist forum, and Hartung et al. (2017) develop a method for identifying right-wing extremist Twitter profiles. However, there is a dearth of NLP research on the more general process of radicalisation. Yet relevant resources exist: recent studies in political science (Baele et al., 2023) and cybersecurity (Vu et al., 2021; Ribeiro et al., 2021) have developed large datasets on online extremism. They address the strongly developed in-group language and imagery using surface features such as the lexicon developed by Farrell et al. (2019).

A challenge is that the concept of “radicalisation” is poorly defined (Della Porta and LaFree, 2012; Schmid, 2016), although it is generally agreed that it involves a gradual process, rather than an instantaneous conversion (Munn, 2019; Bowman-Grieve, 2010). Computational works in this area have tended to treat it as a binary state (eg. Ferrara et al., 2016; Magdy et al., 2016), which ignores this nuance. The lack of a clear definition of the phenomenon further means that human annotation is likely to provide an imperfect and subjective interpretation of the data. Fernandez et al. (2018) have proposed a different approach: looking to behaviour (in particular, the use of terms from an extremist lexicon) as an indicator for how much radical influence an individual is under. This avoids the potentially biased human annotation step, as well as recognising that radicalisation exists along a spectrum. We follow a similar approach in this work, with three further contributions:

- We propose a more holistic approach, considering three dimensions of behaviour: hostile language usage, long-term engagement on an extremist platform, and connectedness within the social network.
- We apply and evaluate modern NLP language modelling techniques, as opposed to count-based methods favoured in prior work.
- We investigate dedicated extremist platforms. Prior work has predominantly focused on Twitter data. Extremist forums are in general operationally different from Twitter, notably lacking a follower graph and user profiles, which necessitates specialised systems.

We proceed by providing a theoretical grounding (Section 2) and formal definition (Section 3) for the three indicators. We further investigate the interaction and development of these factors within anti-women communities (Section 4), which illustrates

that they provide complementary and compelling perspectives. Finally, we investigate the early signs of these indicators, in particular predicting the maximum degree of hostility, longevity and inter-group connectivity measured over an individual’s lifetime, after observing an initial subset of their interactions within the group (Sections 5 and 6).

Our results indicate that it is possible to prioritise at-risk users with a concordance index of 0.70 after 10 posts and 0.68 after 5 posts. Our top-performing approach is a multitask model that jointly predicts the three factors based on a combination of interaction and linguistic inputs. We further investigate the effect of the number of input posts on prediction accuracy, finding a good tradeoff between early prediction and performance is achieved after 6 posts.

## 2 Radicalisation in online communities

In this work, we follow the definition of [Dalgaard-Nielsen \(2010\)](#) as “a **process** in which **radical ideas** are accompanied by the development of a willingness to directly **support** or engage in **violent acts**”, and we specifically focus on radicalisation within online extremist communities.

[Bowman-Grieve \(2010\)](#) argues that the internet can play a role in facilitating individual radicalisation by providing **connection to communities** that reaffirm and strengthen extreme beliefs. They state that members of these communities tend to inhibit various stages in the radicalisation process, and that the formation of interpersonal bonds with radicalised members is an important factor for successful recruitment. According to [Winter et al. \(2020\)](#), linguistic and semantic analysis of online content have been shown to have great potential as part of intelligence-gathering measures; however, they also note that studies in this area have not attempted to identify a definitive set of signals for the potential presence of radicalisation.

The goal of this work is to identify such signals within the scope of online extremist communities. Following the above descriptions, we identify three observable behaviours that relate to online radicalisation at the individual level:

1. Using hostile language originating from a violent extremist ideology (exhibiting adoption of **radical ideas** and **support of violent acts**),
2. Connecting to a network that espouses these extreme ideas (exhibiting **connection to the community**), and

3. A sustained engagement with its doctrine over time (following a **process**).

Existing research has investigated some of these signals in isolation. Targeted hate speech has been used to identify the promoters of various extremist ideologies ([Hartung et al., 2017](#); [Vidgen and Yasseri, 2020](#); [Alatawi et al., 2021](#)). Community connectedness, as measured through network features, has been used to identify key members of terrorist organisations ([Gialampoukidis et al., 2017](#); [Berzinji et al., 2012](#)). In research on communities more broadly, connectedness in the social graph and the adoption of in-group language have been found to be indicative of a user’s likelihood to churn ([Rowe, 2013](#); [Danescu-Niculescu-Mizil et al., 2013](#)), as well as a user’s loyalty to a particular online community ([Hamilton et al., 2017](#)).

A lesser-studied component is the effect of sustained engagement in an extremist group. [Bowman-Grieve \(2010\)](#) states that a sense of status is associated with long-term membership in online extremist communities, and that increased involvement over time may parallel increased ideological development. This notion is also supported by research in psychology: social identity theory holds that group members derive part of their sense of self from the groups to which they belong and will adjust their own behaviours to conform to the group norms ([Hogg and Terry, 2014](#)). Empirical support is provided by [Youngblood \(2020\)](#), who model radicalisation as a social contagion process requiring reinforcement for adoption, and find that social media usage and group membership enhance the spread. [Hassan et al. \(2018\)](#) further find a causal link between membership of Reddit hate groups and the use of hate speech.

Thus, we have identified three radicalisation indicators grounded in prior work: use of hostile language, connectedness in the social graph, and longevity on the platform. In Section 4, we detail how these factors are quantified. Similar to [Fernandez et al. \(2018\)](#) and [Rowe and Saif \(2016\)](#), we do not claim to predict radicalisation, but rather investigate behaviours that may indicate radicalisation. Furthermore, we do not consider these indicators to be exhaustive, but believe that they offer diverse and well-justified perspectives.

## 3 Quantifying radicalisation

We calculate **betweenness centrality** as a measure for the connectedness of an individual in an extrem-

ist community. Betweenness centrality provides a measure of the importance of a node as a function of the number of shortest paths that traverse it, and is often used to identify prominent members of a community (Brandes, 2001). We construct an interaction graph where each node represents a user, and an undirected edge is added between user nodes if they engage in the same conversation thread. The edges are weighted by the number of shared threads. To account for the dynamic nature of the user base, we construct the graph at monthly increments for each community and recalculate the centrality scores for each user. Similar snapshot-based approaches are followed by Hamilton et al. (2017) and Danescu-Niculescu-Mizil et al. (2013). An objection to this approach may be that the coarseness of aggregation might not capture rapid changes in the network; however, it ensures that our models are not overly sensitive to minor fluctuations.

To calculate **hostility**, we use a lexicon of in-group language associated with the community. Extremist factions commonly define themselves through the deliberate exclusion of a specific out-group, and consequently, their internal jargon tends to be hostile towards this out-group. An alternative approach could be to consider a broader definition of hostility using pre-trained toxicity models. However, as mentioned in Section 1, these groups have a propensity for using non-standard in-group language which would not be captured by generalised toxicity models. Lexicon-based approaches are similarly used to investigate radicalisation in Fernandez et al. (2018) and Lara-Cabrera et al. (2017).

**Longevity** is calculated as the number of posts produced by a user in their time on the platform. Time on the platform, in days or months, would also be a possible indicator for longevity and is generally correlated with the volume of posts. However, the former is considered to be a more robust measure as it penalises intermittent and sporadic engagement. A similar argument was adopted by Danescu-Niculescu-Mizil et al. (2013) and Rowe (2013), who quantify the lifecycle stage of users based on the elapsed proportion of their total lifetime post volume, rather than clock time.

## 4 Analysis

In this section, we investigate the indicators described in Section 2 using a dataset of discussions

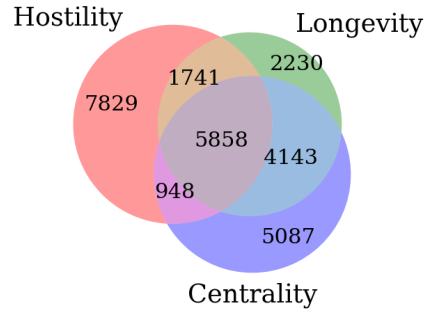


Figure 1: The intersection of the 90th percentile users of longevity, hostility and centrality, showing the number of users per section.

on 8 extremist anti-women forums<sup>1</sup> by Ribeiro et al. (2021). The dataset consists of 7.4 million posts by 139 090 users ranging from 2005 to 2019. For each post, the author, date, thread ID and text are provided. Ribeiro et al. (2021) used this data to study the evolution of different communities over time, whereas this work focuses on the trajectories of individuals.

The forums in this dataset belong to a larger network of online communities collectively referred to as the “manosphere”, which is characterised by sexual objectification of women or endorsements of violence against women. Farrell et al. (2019) and Baele et al. (2023) showed that the language used in manosphere communities is becoming increasingly extreme in nature, and at least 15 acts of real-world terrorism have been connected to this network (Latimore and Coyne, 2023). To measure hostility within this community, we use the lexicon developed by Farrell et al. (2019), consisting of 424 words and phrases. Evaluating the radicalisation indicators on this dataset, a number of conclusions can be drawn.

*(i) Longevity, hostility and centrality provide complementary perspectives.* Figure 1 illustrates the intersection of the 90th percentile users per indicator. To find these groups, we use the maximum indicator value over each user’s lifetime (hereafter referred to as their **eventual** value) and we calculate percentiles for each forum separately. It is evident that the sets intersect to some degree, but there is also substantial non-overlapping components. We further calculate the Spearman correlation between these factors for the full population. The strongest

<sup>1</sup>The dataset also contains posts from anti-women subreddits; however, we chose to focus on single-community dedicated extremist platforms.

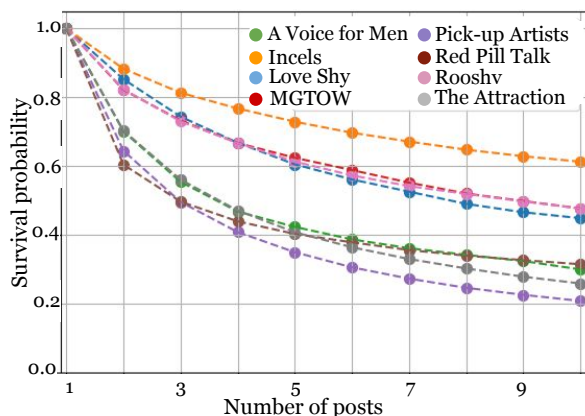


Figure 2: Survival curves for 8 manosphere forums, illustrating the likelihood of a user to continue interacting on the platform after  $N$  posts, for  $N < 10$ .

correlation ( $\rho = 0.798$ ) is observed between the eventual longevity and centrality values, whereas the weakest correlation is between hostility and centrality ( $\rho = 0.469$ ), and  $\rho = 0.613$  for hostility and longevity. All three correlations are statistically significant ( $P \ll 0.05$ ). Thus, we conclude that these factors interact but that each offers a distinct perspective, with hostility being the most disjunct.

*(ii) Many users churn quickly.* There is a steep drop-off in users after relatively few interactions, which aligns with the proposition by Barrelle (2010) that high turnover is characteristic of extreme groups. Figure 2 shows the survival function (Goel et al., 2010) for the number of posts per user for each forum, which illustrates the fraction of users who have more than  $N$  posts, for  $N \leq 10$ . For half of the forums, more than 60% of their users have less than 5 posts in their lifetime. This may be due to users realising after further exposure to the community that the extremeness of the ideology does not resonate with them. The forum with the least churn is Incels, which could be related to the fact that many users migrated to this forum after the *r/incels* subreddit was banned in 2017 (Hauser, 2017); as such, users would already have been inducted into the ideology before joining.

*(iii) Some users start out hostile; others become hostile.* The radicalisation factors vary over the course of a user’s lifetime on the platform. From the positive correlation between hostility and longevity, we know that that users who are on the platform for longer reach higher levels of hostility, but how quickly does this happen? Figure 3 shows the number of days it takes for users to reach the

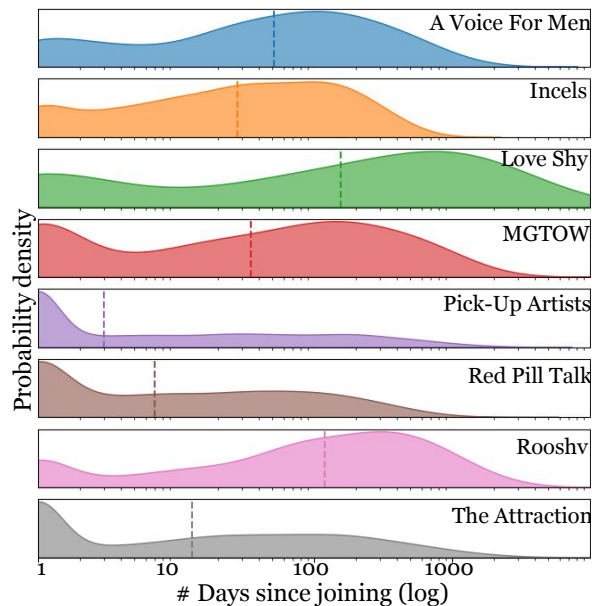


Figure 3: The number of days (logscale) for users to reach the 90th percentile of hostility, per forum.

90th percentile of hostility. For five of the forums, a bimodal distribution is observed, with an early peak ( $< 10$  days) as well as a later peak between 100 and 1000 days. This indicates that a subset of users already exhibit these behaviours when they join the platform, whereas others develop them over time. The stage in their radicalisation process at which a user joins the platform would likely play a role in this phenomenon. This supports the social science research that states that there is no single, agreed upon pathway to radicalisation (Schmid, 2016; Munn, 2019), and highlights the importance of considering multiple indicators.

The three platforms that do not exhibit this trend, having only an early peak, also had higher early churn rates (Figure 2). For the longevity and centrality factors, this bimodality is not present: only a later peak (100–1000 days) is observed.

*(iv) Early signals of eventual behaviour.* Having noted that the indicator values vary over time, we turn to the question of which early signals are predictive of eventual behaviour along the three dimensions. We calculate the following features for the first 10 user interactions for users with 10 or more posts:

- **Post length:** median character count per post,
- **Number of hostility terms:** the median number of terms from the Farrell et al. (2019) lexicon per post,
- **Number of threads** in which a user engaged,

Feature	Centr.	Host.	Long.
Post length	-0.040	0.545	-0.101
# hostility terms	0.156	0.363	0.070
# threads	0.288	-0.075	0.063
Time between posts	-0.184	-0.014	-0.134
# days engaged	0.470	0.468	0.748

Table 1: The Spearman correlation between features of the first 10 posts by a user and eventual indicator levels.

- **Time between posts:** the median number of hours between posts, and
- **Days engaged:** number of distinct days on which the user engaged on the platform.

We calculate the Spearman correlation of the eventual indicator values with the above feature values after 10 interactions. The results, in Table 1, show that these early behaviours are correlated to varying degrees with each of the indicators. All correlations are significant at the  $\alpha = 0.05$  level. A strong correlation to all three indicators is given by the number of distinct days a user engaged on the platform through their first 10 posts. A possible explanation is that a user who comes back repeatedly on separate occasions indicates a higher level of interest and receptiveness, compared to one who posts a larger volume of posts at once, and then disconnects for several days. The largest correlation is to eventual longevity, which aligns with our expectation that longevity is tied to loyalty (Hamilton et al., 2017). Linguistic features (post length and hostility terms) are correlated to eventual hostility, but have no strong relationships to eventual centrality or longevity. Similarly, the number of threads in which a user engaged has a positive correlation to eventual centrality, but a weak relation to longevity and hostility (in a negative direction). This shows that there are early signs of each of the three indicators that are not correlated to the others, providing further support for our multi-indicator approach. The time between posts has a slight negative correlation to centrality and longevity, meaning that more frequent engagements are positively correlated to these indicators.

These results illustrate that there are early signals that preempt users’ eventual behaviour. In the remainder of this paper, we investigate how accurately the three indicators can be predicted.

## 5 Early prediction of indicators

We define the task of predicting a user’s maximum lifetime score on the three radicalisation indicators

after observing an initial subset of  $N$  posts by that user, with  $N \in \{5, 10\}$ . We choose these values of  $N$  based on the survival curves (Fig. 2), which indicate a substantial drop-off in users with less than 5 posts and a stabilisation after  $N = 10$ . Earlier detection is better, but models do require sufficiently strong signals which may not be present if the information is too limited. Since these indicators take on real-valued numbers, this is a regression task.

### 5.1 Metrics

We use two metrics to compare performance on this task. Since an aim of this work is to prioritise users for deradicalisation initiatives, the ordering of users is of interest. To measure this, we report the **concordance index (CI)** (Harrell et al., 1982). A pair of observations  $i, j$  is considered concordant if the prediction and the ground truth have the same inequality relation, i.e.  $(y_i > y_j, \hat{y}_i > \hat{y}_j)$  or  $(y_i < y_j, \hat{y}_i < \hat{y}_j)$ . The concordance index is the fraction of concordant pairs in the test set. A random model would achieve a CI of 0.5 and a perfect score is 1. We also report the mean absolute error (**MAE**) for each indicator. MAE is widely used in regression studies as it provides an intuitive measure for numerical accuracy. However, it is susceptible to outliers and could not be compared between factors, since they operate on different numeric scales. Consequently, we rely on the CI for model selection. Significance testing is performed with the two-sided randomised permutation test, using Monte Carlo approximation with  $R = 9999$ .

### 5.2 Data

We use the Ribeiro et al. (2021) manosphere dataset, described in Section 4, in this evaluation. We filter entries with missing dates, texts, authors or thread IDs and remove users with less than 10 interactions. The resulting dataset contains 7.1 million posts by 39 765 users. The median post length is 33 tokens and the median number of posts per user is 30. The labels are given by the indicator definitions as provided in Section 4 and we release our labels to the community<sup>2</sup>. Since the distributions are heavy-tailed, we truncate the indicator values beyond the 95th percentile of each indicator per forum. We split the data into a training, test and development set with a ratio of 75:15:10.

<sup>2</sup><https://github.com/christinedekock11/radicalisation-indicators>

### 5.3 Methods

Our objective in these experiments is to develop quantitative methods for the early prediction of radicalisation indicators. We therefore experiment with various input and auxiliary task combinations to evaluate their efficacy.

**Feature-based models** We use the features described in Section 4 as a baseline, evaluating models with and without glossary features to investigate the effect of adding linguistic information. For the glossary features, we use the mean and maximum of number of glossary terms per post. The feature and indicator values are normalised using min-max scaling. The model architecture consists of a multi-layer perceptron (MLP) with two hidden layers. Three separate models are trained to predict each indicator value independently. Hyperparameters and training details are provided in Appendix A.

**Text-based models** Models that operate directly upon text, as opposed to engineered features, are expected to capture more nuanced features that extend beyond the hostility lexicon and post length. We use the pretrained `all-mpnet-base-v2`<sup>3</sup> sentence transformer (Reimers and Gurevych, 2019) to obtain an embedding of length 768 for each post. The model architecture consists of an LSTM layer (Hochreiter and Schmidhuber, 1997) followed by two hidden layers. Since the embeddings are produced by a large pretrained language model, we expect that a relatively small number of layers should be sufficient to finetune them to our task.

**Mixed-input models** A dual-input architecture is used to combine the text-level learning from embeddings with the engineered interaction and glossary-based features. The glossary-based features capture the use of non-standard in-group terms which may not appear in the vocabulary of a pretrained language model; as such, both types of linguistic inputs may be useful. An LSTM layer and two MLP layers are used to process the text and feature inputs in parallel. The outputs are concatenated and two further hidden layers are applied.

**Multitask models** The analysis in Section 4 indicated that the different indicators interact and correlate to some extent. As such, we expect that parameter sharing might be beneficial, as opposed to training a separate model for each indicator. We

keep the same initial architecture as in the mixed input models, but use a separate prediction head with two additional hidden layers for each output.

Our dataset consists of user profiles from 8 platforms, which may have distinct user-level characteristics. To investigate whether there are useful features that are tied to the different platforms, we further experiment with predicting the forum from which the sample originates as an auxiliary task.

**Survival regression** For time-to-event prediction from text inputs, such as the longevity prediction task, survival regression has been illustrated to outperform traditional regression approaches (De Kock and Vlachos, 2021). This framework has a more explicit treatment of time and events within a standard regression setting, and is particularly effective for modelling real-valued, exponentially-distributed outcomes. We use the logistic hazard model (Gensheimer and Narasimhan, 2019) for the longevity predictions. This framework enables us to retain the same neural architectures, but modify the objective to predict the probability of churn for an individual within each timestep, given survival up to that point (also known as the hazard). The outputs are transformed into 100 equidistant timesteps, and the loss is the negative log likelihood of the predicted versus actual hazard per timestep.

## 6 Results

Our results are shown in Table 2. Significance of improvements in CI ( $P \leq 0.05$ ) as compared to the model directly above is indicated by asterisks. The CI scores for the three indicators are in a relatively close range to one another for most models. The top-performing model has a CI of 0.667 for centrality, 0.698 for hostility and 0.681 for longevity (at  $N = 10$ ), constituting a statistically significant improvement over baselines of respectively +1%, +6.3% and +7.9%. For all models and indicators, the performance at  $N = 5$  is worse than at  $N = 10$ . Of the three indicators, centrality has the largest increase in CI between  $N = 5$  and  $N = 10$ . The MAE values generally follow the CIs in terms of direction of improvement.

Adding sources of information or auxiliary tasks tends to improve performance in our experiments. Using glossary-based features in addition to interaction-based features improves CI (significant for 4 out of 6 cases), which supports our central hypothesis that linguistic cues can be helpful at foreshadowing radicalisation. Using only post

<sup>3</sup><https://huggingface.co/sentence-transformers/all-mpnet-base-v2>.

Model	Centrality		Hostility		Longevity	
	CI $\uparrow$	MAE $\downarrow$	CI $\uparrow$	MAE $\downarrow$	CI $\uparrow$	MAE $\downarrow$
$N = 5$						
Interaction features	0.620	0.380	0.616	7.150	0.561	49.43
Interaction + glossary features	0.621	0.388	0.640*	7.258	0.572*	50.46
Transformer embeddings	0.595	0.376	0.658*	7.628	0.647*	46.33
+ all features	0.608*	0.381	0.666	7.754	0.652	46.55
+ multifactor training	<b>0.622*</b>	0.315	0.672	5.730	0.645	<b>45.18</b>
+ forum aux. task	0.621	<b>0.314</b>	<b>0.677</b>	<b>5.737</b>	<b>0.656*</b>	45.675
$N = 10$						
Interaction features	0.657	0.388	0.635	7.279	0.602	48.15
Interaction + glossary features	0.659	0.390	0.665*	7.341	0.615*	47.59
Transformer embeddings	0.616	0.382	0.679*	7.749	0.654*	45.12
+ all features	0.651*	0.393	0.689	7.956	0.677*	44.40
+ multifactor training	0.666*	<b>0.287</b>	0.693	<b>5.527</b>	0.672	43.56
+ forum aux. task	<b>0.667</b>	0.288	<b>0.698</b>	5.538	<b>0.681*</b>	<b>43.24</b>

Table 2: Results for predicting the eventual centrality, hostility and longevity values at  $N = 5$  and  $N = 10$ . Arrows indicate the preferred directions per metric and best models per indicator and metric are shown in bold. Significance of improvements in CI ( $P \leq 0.05$ ) as compared to the model directly above is indicated by asterisks.

embeddings outperforms feature-based approaches for hostility and longevity prediction, but reduces the CI for centrality. Combining features and embeddings improves the CI over embedding-only models (significant for 3 out of 6 cases), indicating that the features contain useful information beyond what is captured by the language model. Joint training of the three indicators yields a further improvement, particularly in MAE, which aligns with expectation that the three factors contain mutually informative signals. Marginal improvements, significant in 2 cases, are made by adding the forum prediction auxiliary task. The experiments in the remainder of this section use this model.

The performance of the feature-based centrality model declined when the text embeddings were added, and although the highest score for this indicator was achieved by the multifactor model which uses embeddings, this improvement was smaller than for the other indicators. Considering that the analysis in Table 1 showed no correlation between the early use of hostility terms and eventual centrality, this is perhaps not surprising. We can conclude that the language features and models used in this study are less apt at detecting the early cues that foreshadow centrality, if they are present.

## 6.1 Optimising the number of inputs

Our aim in this work is the early identification of users who are at risk of radicalisation. In this section, we consider *how early* such a prediction might be made. Given the tradeoff between prioritising performance versus earlier prediction, the optimal prediction point will be where improvement starts to saturate as  $N$  increases. To find this,

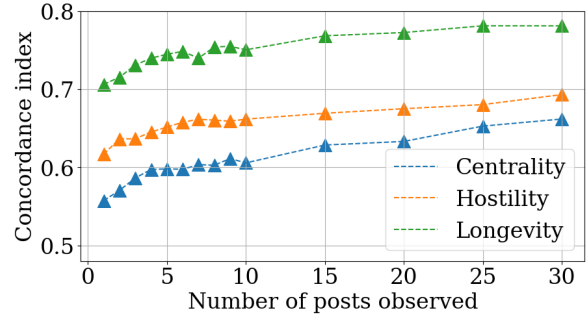


Figure 4: Performance at different  $N$ .

	2	3	4	5	6	7	8	9	10
1	<b>.029</b>	.037	0	0	0	0	0	0	0
2	-	.994	.325	.093	<b>.02</b>	.004	.01	.013	.009
3	-	-	.323	.098	.017	<b>.006</b>	<b>.005</b>	<b>.007</b>	<b>.006</b>
4	-	-	-	.475	.167	.078	.105	.119	.082
5	-	-	-	-	.47	.276	.316	.419	.268
6	-	-	-	-	-	.65	.755	.86	.652
7	-	-	-	-	-	-	.907	.791	.988
8	-	-	-	-	-	-	-	.88	.875
9	-	-	-	-	-	-	-	-	.767

Table 3: Significance of performance increases with larger  $N$  for the hostility indicator.

we train models with inputs ranging from 1 to 30 posts, sampling more densely at  $N < 10$  as larger improvements are expected.

The results are shown in Figure 4. Only users with 30 or more posts are included in this experiment, so the CI values cannot be directly compared to the results in Table 2. For all three indicators, there is an upward trend in CI as  $N$  increases, with a steeper increase for  $N < 5$  and a more moderate improvement for  $5 < N \leq 10$ . Beyond  $N = 10$ , diminishing returns are observed for the longevity and hostility indicators, meaning that delaying the

Training data	Manosphere			Stormfront		
	Cent	Host	Long	Cent	Host	Long
Manosphere	0.666	0.693	0.672	0.592	0.660	0.584
Stormfront	–	–	–	0.635*	0.682*	<b>0.603*</b>
Combined	0.662	0.689	0.667	0.635	0.705*	0.590
+ forum task	<b>0.668</b>	<b>0.699</b>	<b>0.675</b>	<b>0.640*</b>	<b>0.721*</b>	0.598

Table 4: Concordance index of multifactor models for the Manosphere and Stormfront datasets.

prediction beyond this point is not well-justified. It is worth noting that centrality still improves substantially beyond this point.

We are interested in the minimum improvement in  $N$  which would constitute a significant improvement in CI. We use randomised permutation testing to evaluate the significance of the improvement at each step for  $N < 10$ . The P-values for hostility are shown in Table 3, with significance ( $P \leq 0.05$ ) indicated in green. A significant improvement ( $P = 0.029$ , shown in bold) is observed between 1 and 2 inputs. From 2, we would need to increase the number of inputs to 6 to obtain a significant improvement ( $P = 0.02$ ). No further significant improvements are possible in the observed range. For centrality and longevity, following a similar procedure yields significant improvements until  $N = 8$  and  $N = 6$ , respectively. As such, we recommend using the initial 6 posts made by a user to predict radicalisation as early as possible with a good tradeoff in accuracy.

## 6.2 Application in other communities

This paper is concerned with radicalisation as a general concept, and not only its specific manifestation in the manosphere. As such, we also evaluate our framework on the white supremacy platform Stormfront, using the ExtremeBB dataset (Vu et al., 2021). Applying the same filters as in Section 5.2, we obtain a dataset of posts by 25 895 users. The centrality and longevity indicators are calculated as described in Section 3. The hostility indicator is intended to capture the adoption of extreme ideas from the community in question, which we operationalise using a lexicon. A list of 293 alt-right phrases and symbols was scraped from Rational-Wiki<sup>4</sup> and is shared with the community. The indicator labels for this dataset cannot be shared under the ExtremeBB data agreement.

We expect to see differences in the numeric values of the indicators as their distributions will differ

<sup>4</sup>[https://rationalwiki.org/wiki/Alt-right\\_glossary](https://rationalwiki.org/wiki/Alt-right_glossary)

between the populations. This is accounted for in our framework by (i) applying min-max scaling to the indicator values during training, and (ii) using the CI metric for evaluation, which is concerned with relative ordering rather than absolute values.

We evaluate a number of different training configurations, with CI values at  $N = 10$  shown in Table 4. Using the best model as trained on manosphere data, lower CI values are recorded for all three indicators compared to the original dataset. Training on the Stormfront dataset instead improves the scores for all three indicators on the same data (significant at the  $\alpha = 0.05$  level). Training on both datasets increases the CI for the hostility prediction on Stormfront but reduces the CI for all others. However, when the forum prediction auxiliary task is included, there is a statistically significant improvement on the centrality and hostility metrics on the Stormfront data.

In conclusion, a drop in model performance is to be expected if a model trained on data from one extremist community is transferred to a different community without any adjustment. However, joint training on unrelated communities is useful if the platform information is provided in the form of an auxiliary task. Future work may explore training on larger multi-community datasets.

## 7 Conclusion

We have proposed a framework for quantifying behaviours that are indicative of radicalisation. We investigated the interaction of these indicators using a dataset of posts on extremist platforms and identified early signals that correspond to the eventual indicator levels of an individual. We then developed and evaluated models that can preemptively rank potentially at-risk users.

A comprehensive understanding of radicalisation requires inputs from several disciplines to capture the various contributing factors, including the psychological, educational, economic, and social-adjustment parameters of the individual. Capturing these factors in a single predictive model is not feasible within the current data landscape. Using behaviour as a proxy for some of these parameters, identifying the most predictive attributes, and modelling them using NLP is a promising methodology. We look forward to addressing more of these parameters in work across relevant disciplines.



## 8 Limitations

We hope that this work will serve as a foundation for further NLP work in this direction, which may address some of the following limitations.

The hostility indicator is reliant on a lexicon, which is a standard practice for work in this space. Linguistic resources have been developed for many online extremist communities. However, using manually constructed lexicons is sub-optimal as they are bound to have imperfect recall and they are constructed for the community at a particular point in time, which ignores the fact that community language is highly dynamic.

The centrality indicator is intended to capture social connectedness and is a well-established metric for this purpose. However, extremist groups are known to be prone *splintering*, a process whereby the more extreme community members form sub-groups with limited interaction with the larger community. This behaviour is highly indicative of radicalisation but is not captured by the centrality indicator.

The longevity metric assumes that users who churn early, do so because they are disengaging from the group. It is also plausible that some users may leave a community to seek out more extreme groups. However, since early churn is commonly observed in all extreme groups (Barrelle, 2010), we assume that the former explanation holds true for the majority of users.

Finally, our work builds on prior research in online communities. More consideration could be devoted to the characteristics that differentiate extreme communities from online communities more broadly.

## 9 Ethics

A motivation of our work is the ability to monitor discussions and identify at-risk users in online extremist communities. It could conceivably be misused to profile and pre-emptively prosecute individuals. Since our evaluation shows that the predictive models are not perfectly accurate, that would be a gross abuse of the technology, and we do not release our models publicly to mitigate this risk. However, the models can be useful as a part of larger intelligence gathering systems, as mentioned by Winter et al. (2020).

We would further like to reiterate that these are not general purpose approaches for online discussions, and that the indicators would not make sense

to signify radicalisation within more general social networks, where people engage on various topics. We are specifically looking at individuals in dedicated extremist forums, and aiming to anticipate how much they will become entrenched in the community and express ideas from the extremist ideology.

## References

- Hind S. Alatawi, Areej M. Alhothali, and Kawthar M. Moria. 2021. [Detecting white supremacist hate speech using domain specific word embedding with deep learning and bert](#). *IEEE Access*, 9:106363–106374.
- Stephane Baele, Lewys Brace, and Debbie Ging. 2023. A diachronic cross-platforms analysis of violent extremist language in the incel online ecosystem. *Terrorism and Political Violence*, pages 1–24.
- Kate Barrelle. 2010. Disengagement from violent extremism. In *Conference paper. Monash University: Global Terrorism Research Centre and Politics Department*.
- Ala Berzinji, Lisa Kaati, and Ahmed Rezine. 2012. Detecting key players in terrorist networks. In *2012 European Intelligence and Security Informatics Conference*, pages 297–302. IEEE.
- Khandis R Blake, Siobhan M O’Dean, James Lian, and Thomas F Denson. 2021. Misogynistic tweets correlate with violence against women. *Psychological science*, 32(3):315–325.
- Lorraine Bowman-Grieve. 2010. The internet and terrorism: pathways towards terrorism & counter-terrorism. In Andrew Silke, editor, *The psychology of counter-terrorism*. Routledge.
- Ulrik Brandes. 2001. A faster algorithm for betweenness centrality. *Journal of mathematical sociology*, 25(2):163–177.
- Anja Dalgaard-Nielsen. 2010. Violent radicalization in europe: What we know and what we do not know. *Studies in conflict & terrorism*, 33(9):797–814.
- Cristian Danescu-Niculescu-Mizil, Robert West, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013. No country for old members: User lifecycle and linguistic change in online communities. In *Proceedings of the 22nd international conference on World Wide Web*, pages 307–318.
- Ona de Gibert, Naiara Perez, Aitor García-Pablos, and Montse Cuadros. 2018. [Hate speech dataset from a white supremacy forum](#). In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 11–20, Brussels, Belgium. Association for Computational Linguistics.

- Christine De Kock and Andreas Vlachos. 2021. [Survival text regression for time-to-event prediction in conversations](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1219–1229, Online. Association for Computational Linguistics.
- Donatella Della Porta and Gary LaFree. 2012. Guest editorial: Processes of radicalization and de-radicalization. *International Journal of Conflict and Violence (IJCV)*, 6(1):4–10.
- Tracie Farrell, Miriam Fernandez, Jakub Novotny, and Harith Alani. 2019. Exploring misogyny across the manosphere in reddit. In *Proceedings of the 10th ACM conference on web science*, pages 87–96.
- Miriam Fernandez, Moizzah Asif, and Harith Alani. 2018. Understanding the roots of radicalisation on twitter. In *Proceedings of the 10th ACM conference on web science*, pages 1–10.
- Emilio Ferrara, Wen-Qiang Wang, Onur Varol, Alessandro Flammini, and Aram Galstyan. 2016. Predicting online extremism, content adopters, and interaction reciprocity. In *Social Informatics: 8th International Conference, SocInfo 2016, Bellevue, WA, USA, November 11-14, 2016, Proceedings, Part II 8*, pages 22–39. Springer.
- Michael F Gensheimer and Balasubramanian Narasimhan. 2019. A scalable discrete-time survival model for neural networks. *PeerJ*, 7:e6257.
- Ilias Gialampoukidis, George Kalpakis, Theodora Tsirikla, Symeon Papadopoulos, Stefanos Vrochidis, and Ioannis Kompatsiaris. 2017. [Detection of terrorism-related twitter communities using centrality scores](#). In *Proceedings of the 2nd International Workshop on Multimedia Forensics and Security, MFSec '17*, page 21–25, New York, NY, USA. Association for Computing Machinery.
- Paul Gill, Emily Corner, Maura Conway, Amy Thornton, Mia Bloom, and John Horgan. 2017. Terrorist use of the internet by the numbers: Quantifying behaviors, patterns, and processes. *Criminology & Public Policy*, 16(1):99–117.
- Manish Kumar Goel, Pardeep Khanna, and Jugal Kishore. 2010. Understanding survival analysis: Kaplan-meier estimate. *International journal of Ayurveda research*, 1(4):274.
- William Hamilton, Justine Zhang, Cristian Danescu-Niculescu-Mizil, Dan Jurafsky, and Jure Leskovec. 2017. Loyalty in online communities. In *Proceedings of the International AAAI conference on web and social media*, volume 11, pages 540–543.
- Frank E Harrell, Robert M Califf, David B Pryor, Kerry L Lee, and Robert A Rosati. 1982. Evaluating the yield of medical tests. *Jama*, 247(18):2543–2546.
- Matthias Hartung, Roman Klinger, Franziska Schmidtke, and Lars Vogel. 2017. [Ranking right-wing extremist social media profiles by similarity to democratic and extremist groups](#). In *Proceedings of the 8th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 24–33, Copenhagen, Denmark. Association for Computational Linguistics.
- Ghayda Hassan, Sébastien Brouillette-Alarie, Séraphin Alava, Divina Frau-Meigs, Lysiane Lavoie, Arber Fetiu, Wynnnpaul Varela, Evgueni Borokhovski, Vivek Venkatesh, Cécile Rousseau, et al. 2018. Exposure to extremist online content could lead to violent radicalization: A systematic review of empirical evidence. *International journal of developmental science*, 12(1-2):71–88.
- Christine Hauser. 2017. Reddit bans ‘incel’ group for inciting violence against women. *New York Times*. Accessed 10-11-2023.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Michael A Hogg and Deborah J Terry. 2014. *Social identity processes in organizational contexts*. Psychology Press.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Raúl Lara-Cabrera, Antonio Gonzalez-Pardo, and David Camacho. 2017. Statistical analysis of risk assessment factors and metrics to evaluate radicalisation in twitter. *Future Generation Computer Systems*, 93:971–978.
- Jasmine Latimore and John Coyne. 2023. Incels in australia: the ideology, the threat, and a way forward.
- Walid Magdy, Kareem Darwish, and Ingmar Weber. 2016. Failedrevolutions: Using twitter to study the antecedents of isis support. *First Monday*, 21(2).
- Luke Munn. 2019. Alt-right pipeline: Individual journeys to extremism online. *First Monday*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992.
- Manoel Horta Ribeiro, Jeremy Blackburn, Barry Bradlyn, Emiliano De Cristofaro, Gianluca Stringhini, Summer Long, Stephanie Greenberg, and Savvas Zannettou. 2021. The evolution of the manosphere across the web. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 15, pages 196–207.

Matthew Rowe. 2013. [Mining user lifecycles from online community platforms and their application to churn prediction](#). In *2013 IEEE 13th International Conference on Data Mining*, pages 637–646.

Matthew Rowe and Hassan Saif. 2016. Mining pro-  
pensity radicalisation signals from social media users. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 10, pages 329–338.

Alex P Schmid. 2016. Research on radicalisation: Top-  
ics and themes. *Perspectives on terrorism*, 10(3):26–  
32.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky,  
Ilya Sutskever, and Ruslan Salakhutdinov. 2014.  
Dropout: a simple way to prevent neural networks  
from overfitting. *The journal of machine learning  
research*, 15(1):1929–1958.

Bertie Vidgen and Taha Yasseri. 2020. Detecting weak  
and strong islamophobic hate speech on social me-  
dia. *Journal of Information Technology & Politics*,  
17(1):66–78.

Anh V Vu, Lydia Wilson, Yi Ting Chua, Ilia Shumailov,  
and Ross Anderson. 2021. Extremebb: Enabling  
large-scale research into extremism, the manosphere  
and their correlation by online forum data. *arXiv  
preprint arXiv:2111.04479*.

Charlie Winter, Peter Neumann, Alexander Meleagrou-  
Hitchens, Magnus Ranstorp, Lorenzo Vidino, and  
Johanna Fürst. 2020. Online extremism: re-  
search trends in internet activism, radicalization, and  
counter-strategies. *International Journal of Conflict  
and Violence (IJCV)*, 14:1–20.

Mason Youngblood. 2020. Extremist ideology as a com-  
plex contagion: the spread of far-right radicalization  
in the united states between 2005 and 2017. *Humanities  
and Social Sciences Communications*, 7(1):1–10.

## A Training specifications

In all experiments, we use a batch size of 32 and  
ReLU activation functions between hidden layers.  
We train with early stopping with a patience of 20  
epochs. Models are developed in PyTorch. We use  
a gridsearch to determine the best hyperparameter  
values, experimenting with hidden layer sizes in  
{32, 64, 128} and dropout (Srivastava et al., 2014)  
with  $p \in \{0.1, 0.2, 0.5\}$ . The Adam (Kingma and  
Ba, 2014) optimiser is used, with  $\eta \in \{1e-4, 5e-  
4, 1e-3\}$ . The best value per model are reported  
in Tables 5.

<b>Model</b>	<b>Factor</b>	<b>Dropout (<math>p</math>)</b>	<b>Hidden units per layer</b>	<b>Learning rate</b>
$N = 5$				
Frequency features	Centrality	0.1	32	0.0005
	Hostility	0.2	32	0.0005
	Longevity	0.2	128	0.0005
Frequency + glossary features	Centrality	0.1	64	0.0005
	Hostility	0.1	32	0.0005
	Longevity	0.1	64	0.0005
Embeddings	Centrality	0.1	32	0.0001
	Hostility	0.1	64	0.0001
	Longevity	0.2	128	0.0005
Embeddings + features	Centrality	0.1	64	0.0005
	Hostility	0.1	32	0.0001
	Longevity	0.1	64	0.0005
Multifactor + forum aux.task	All	0.1	64	0.0005
	All	0.1	128	0.0005
$N = 10$				
Frequency features	Centrality	0.1	128	0.0005
	Hostility	0.1	32	0.0005
	Longevity	0.2	128	0.0005
Frequency + glossary features	Centrality	0.1	32	0.0005
	Hostility	0.1	128	0.0005
	Longevity	0.1	32	0.0005
Embeddings	Centrality	0.1	32	0.0001
	Hostility	0.2	64	0.0001
	Longevity	0.1	128	0.0005
Embeddings + features	Centrality	0.1	32	0.0001
	Hostility	0.2	64	0.0001
	Longevity	0.1	128	0.0005
Multifactor + forum aux.task	All	0.1	128	0.0005
	All	0.1	128	0.0001

Table 5: Hyperparameters for per-factor models.