# POSTCOG: A Tool for Interdisciplinary Research into Underground Forums at Scale

Ildiko Pete, Jack Hughes, Andrew Caines, Anh V. Vu, Harshad Gupta,
Alice Hutchings, Ross Anderson, Paula Buttery
*firstname.lastname@cl.cam.ac.uk*
Department of Computer Science & Technology
University of Cambridge, Cambridge, CB3 0FD, UK

*Abstract*—Underground forums provide useful insights into cybercrime, where researchers analyse underlying economies, key actors, their discussions and interactions, as well as different types of cybercrime. This interdisciplinary topic of study incorporates expertise from diverse areas, including computer science, criminology, economics, psychology, and other social sciences. Historically, there were significant challenges around access to data, but there are now research datasets of millions of messages scraped from underground forums. The problems now stem from the size of these datasets and the technical nature of methods and tools available for data sampling and analysis at scale, which make data exploration difficult for non-technical users. POSTCOG has been developed to solve this problem. We first provide a survey of prior work into underground forums; this was used to understand the requirements and functionalities valued by researchers, and to inform the design of a data exploration tool. We then describe POSTCOG, a web application developed to support users from both technical and non-technical backgrounds in forum analyses, such as search, information extraction and cross-forum comparison. The prototype's usability is then evaluated through two user studies with expert users of the CRIMEBB dataset. POSTCOG is made available for academic research upon signing an agreement with the Cambridge Cybercrime Centre.

## I. INTRODUCTION

Underground forums play a central role in cybercrime and related online communities. They provide a platform for people engaging in illicit activities to find and interact with each other, to share ideas and mutual interests or know-how. For example, members of HACK FORUMS, one of the largest English language underground communities, have shared more than 61 million posts since 2007. While a number of HACK FORUMS members have been prosecuted for cybercrime offences [49], one notable forum member, Michael Hutchins, helped stop the spread of WannaCry ransomware in 2017 [34].

As cybercrime is a complex phenomenon, research benefits from an interdisciplinary approach, incorporating methods and theoretical perspectives from criminology, sociology, economics and other social sciences, as well as computer science. Both quantitative and qualitative approaches provide insights at the macro level (across forums, across time) and the micro level (exploring a topic in detail). By combining these methods, we can explore topics such as the pathways of key actors [49], the evolution of cybercrime markets [70], and cybercrime activities relating to the Internet of Things [5].

However, due to the sheer scale of these forums, manual analysis techniques are inadequate. Identifying and extracting relevant data for a specific study can be a huge challenge, particularly for interdisciplinary researchers who are not accustomed to working with large datasets. They urgently need better tools for extracting data of interest at scale, where machine learning (ML), specifically natural language processing (NLP) techniques, are promising approaches [10], [26], [53].

Another problem for underground forum research is the difficulty in obtaining authentic data. To solve this problem and enable large scale longitudinal analysis of underground forums, the Cambridge Cybercrime Centre[1] has been collecting and sharing underground forum datasets (and various other cybercrime and extremism collections) with the research community. The CRIMEBB dataset currently contains data from 34 underground forums in five languages (English, Russian, German, Arabic, and Spanish) [50].

This resource is particularly useful for social science researchers, who might be deterred from collecting their own data by the technical difficulties of web scraping. At the time of writing, CRIMEBB is 121 GB, representing over 99 million posts, and the database continues to grow. However, even with access to the data, nontrivial computational tasks of information retrieval and extraction must be carried out in order to work with targeted subsets of the database. Working with big data presents significant resource challenges from storage and access through the execution of search queries, and these remain a barrier to many researchers.

To this end, we have developed a web application, named POSTCOG, to reduce the barriers and enable underground forum analysis by scholars from a variety of academic disciplines. POSTCOG achieves this goal by (1) providing a web interface for exploring underground forum data and by (2) integrating NLP tools to allow the automated analysis of posts. While many text analysis tools are available, applying them to underground forums requires further customisation due to the language used by members; slang, technical jargon, and abbreviations are commonplace, as well as deliberate obfuscation (e.g. leetspeak). Thus, custom NLP tools built for classifying posts into different types (e.g., [10]) are particularly useful for users who wish to gain insights into these forums.

This paper presents POSTCOG and demonstrates how it supports users in their data exploration and analysis. We start with a literature review of prior work on underground forums to identify researchers' needs in §II, then briefly introduce our CRIMEBB dataset and draw an overview of the data analysis workflow for underground forums in §III. We describe POSTCOG – a toolkit for interdisciplinary analyses of underground forums in §IV, before presenting an evaluation

---

[1] Cambridge Cybercrime Centre: https://www.cambridgecybercrime.uk/

of the tool's prototypes with experienced researchers who are actively using the CRIMEBB dataset in §V.

## II. UNDERGROUND FORUM ANALYSIS AND DATA

Performing data analysis on an underground forum dataset presents a number of challenges to the uninitiated. The existing analysis tools usually involve a steep learning curve [51]. Researchers from non-technical backgrounds may thus be deterred from carrying out the analysis themselves, or may fall back on manual analysis techniques. This limits them to very small subsets, and raises serious questions around sample selection. To understand the landscape and data analysis process in underground forum research, and to highlight areas where researchers from non-technical backgrounds face real barriers, we conducted a brief survey of the field. We include all papers which name the forums captured in the CRIMEBB database that we used to develop POSTCOG.

The survey began by searching for an initial set of research papers using the keywords 'underground forum analysis' and 'dark web forum analysis', with irrelevant results excluded. General analysis work centred around forums such as Twitter and Reddit was discarded, and as was dark-web forum analysis focusing on radicalisation and terrorism. This left us with 64 papers, from which we conducted backward and forward snowballing. We finally pre-selected papers to include only those relating to forums in the CRIMEBB dataset [2]. Some researchers chose not to name the forums analysed in their papers, which meant they are excluded. Our survey is based on a final selection of 60 papers, and aims to explore five main aspects: areas of focus; data sources and volume; data extraction; methods used; and tools and programming languages.

### A. Areas of focus

Researchers are interested in various problems related to underground forums. Some are interested in characteristics of individuals, groups, subcommunities or the entire community; others analyse the contents of posts to understand specific phenomena and forum activities. Here we give a brief description of the main topics.

*1) Identifying key actors:* Identifying and characterising important forum users provides useful insights into the communities in which they operate. Their post contents and activity patterns help us understand the evolution of cyber attacks, the tactics, tools and procedures in use, the incentives facing various actors, and how new actors are recruited and developed. This can support intervention and prevention campaigns [49]. Unsurprisingly, there is considerable interest [49], [27], [76], [77], [25], [31] in key actors and this is a core theme of underground forum analysis.

---

[2] Surface web forums: HACK FORUMS, FREEHACKS, OFFENSIVE COMMUNITY, STRESSER FORUMS, MULTIPLAYER GAME HACKING, LOLZTEAM, GREYSEC, OGUSERS, SAFE SKY HACKS, V3RMILLION, FORUM TEAM, UNKNOWNCHEATS, UNDERC0DE, ZISMO, PROBIV, ANTICHAT, GARAGE4HACKERS, INDETECTABLES, ELHACKER, IFUD, XSS, HACKERS ARMIES, RAIDFORUMS, BLACKHATWORLD, NULLED, CRACKED, KERNELMODE; dark web forums: TORUM, DREAD, DEUTSCHLAND IM DEEP WEB, ENVOY, PIRATEBAY, RUNION, and THE HUB

*2) Analysing underground economies:* Underground forums are not only a platform to exchange ideas and information, but also serve as marketplaces for trading. Investigating their transactions and related discussions can help us understand how these communities operate. Papers in this category focus on analysing markets [70], [53], [64], supply chains [7], currency exchanges and the underground economy generally including the mix of products being sold [35], [17], [18].

*3) Analysing and identifying different types of cybercrime:* Underground forums are used to discuss, promote, and market different types of criminal services. Forum members trade booter services [29], malware, accounts [19], kits for eWhoring [30], and other fraudulent activities [62]. Understanding these is a central problem for underground forum analysis.

*4) Analysing social networks on forums:* Analysing communities within underground forums helps us understand how these communities operate and how members are connected to each other. By modelling these communities as networks we can study them using network topology or other large-scale structural properties, which provide insights into information flow and interaction patterns. There is now a wide range of methods in the toolbox of network science and social network analysis (SNA) [32], [22], [1], [44].

*5) Analysing forum discussions:* A significant number of the studied papers analyse underground forum discussions, which are centred around the following problems:

- Analysis of topics of discussion helps with such tasks as identifying cybercrime-suspect threads [75] and gaining insight into specific attacks [74], [5]. A growing body of work detects trending terms, identifying the topics being discussed on forums, in specific bulletin boards, in public posts or in private messages [26].
- Analysis of language can help us to identify cybersecurity related words or posts for threat intelligence [14], to extract neologisms [37], to carry out authorship attribution and duplicate account detection [2], or to identify dark jargon [59], [73], [60].
- Exploring of the perception of gender, given the role played by misogyny in some crime types and criminal subcultures [4].
- Automatic classification of discussions based on post type and sentiment [10], or for cyber threat intelligence [72].
- Analysis of source code shared on forums [3].

*6) Exploring private messages:* Private messages are studied less than public posts. The challenge is access to data, with researchers typically using leaked or seized datasets. Such analysis requires careful ethical consideration, given that the authors of private messages did not intend them to be publicly available [47], [63].

### B. Data sources and volume

Datasets have been obtained by researchers in a number of ways. Some collect data themselves [75], [77], [25], for example by scraping underground forums. Others request data from other sources, or work with publicly available datasets, such as databases that have been breached and publicly leaked.

Table I: Scope of analysis and methods of research papers included in our survey.

| Research Papers | Extent of Sample | | Data Sources | | Method of Analysis | | | |
|---|---|---|---|---|---|---|---|---|
| | Multiplicity⋆ | Completeness⋄ | Leaked | Scraped | Analytical | SNA | Modelling | Qual |
| [14], [13] | ◑ | ◑ | ✓ | · | · | · | ✓ | · |
| [10], [26], [75] | ◑ | ◑ | · | ✓ | · | · | ✓ | · |
| [30] | ◑ | ◑ | · | ✓ | · | · | · | ✓ |
| [70], [27] | ◑ | ◑ | · | ✓ | ✓ | · | ✓ | · |
| [4] | ◑ | ◑ | · | ✓ | ✓ | · | · | ✓ |
| [40] | ◑ | ◑ | · | ✓ | ✓ | · | · | · |
| [31] | ◑ | ● | ✓ | · | · | ✓ | · | · |
| [18], [32], [63] | ◑ | ● | ✓ | · | · | · | ✓ | · |
| [62] | ◑ | ● | ✓ | · | ✓ | · | ✓ | ✓ |
| [77], [35] | ◑ | ● | · | ✓ | · | · | ✓ | · |
| [49], [46] | ◑ | ● | · | ✓ | · | ✓ | · | · |
| [19] | ● | ◑ | ✓ | ✓ | · | · | ✓ | · |
| [76], [7], [17], [37], [59], [73], [60], [72], [3], [55], [12], [45], [57] | ● | ◑ | · | ✓ | · | · | ✓ | · |
| [47] | ● | ◑ | ✓ | · | · | · | ✓ | · |
| [25] | ● | ◑ | · | ✓ | ✓ | ✓ | · | · |
| [5], [43], [21] | ● | ◑ | · | ✓ | · | · | · | ✓ |
| [53], [74] | ● | ◑ | · | ✓ | ✓ | · | ✓ | · |
| [2] | ● | ◑ | ✓ | · | ✓ | · | ✓ | · |
| [15], [24] | ● | ◑ | · | ✓ | ✓ | · | · | · |
| [52], [41] | ● | ◑ | · | ✓ | · | ✓ | · | · |
| [1] | ● | ● | ✓ | · | ✓ | · | · | · |
| [22] | ● | ● | · | ✓ | · | ✓ | · | · |
| [44] | ● | ● | ✓ | · | · | ✓ | · | · |
| [65], [54], [58] | ● | ● | · | ✓ | · | · | ✓ | · |
| [48] | ● | ● | · | ✓ | · | · | · | ✓ |

⋆**Multiplicity**: ● multiple forums ◑ single forum. ⋄**Completeness**: ● full dataset ◑ sub dataset.

For instance, the NULLED forum database leak is widely used [14], [47], [13], [31].

Another important aspect of working with datasets is the format in which they are available. This determines how steep the learning curve is for researchers from non-technical backgrounds. We found that datasets are mostly made available in a database or SQL dump format [49], [30], [10], [1], [44], [2], which can be used to restore original databases. Depending on the tasks, researchers extract and analyse varying amounts of data, and scraped data can also vary in size. Volume can be expressed in terms of the number of posts, threads and users. The largest number of posts in this analysis is in excess of 32 million posts [15], while the smallest is 840 posts [43].

### C. Data extraction

In the sampling and extraction step, raw data is exported from the original data source, preserving it for further analysis. Depending on the research questions and units of analysis, extracted data may include posts, threads, or forum users. Researchers may use the entire dataset, containing all the activity from one or more forums, a random sample, or only extract specific data relevant to a topic of interest. These data items can take the form of an entire post or specific information, such as currencies, or terms of interest. Some researchers work with the full dataset, but this is rare due to the volume of data.

*1) Keyword-based thread extraction:* Researchers might be interested not only in individual posts but their context. In such cases they extract full threads based on keywords in thread titles or post contents. Appropriate keywords are usually identified by domain experts. Examples of such an approach include the extraction of eWhoring tutorials from thread titles using the keywords '[TUT]' and 'guide' [30], or identifying the types of items available on forums by extracting thread titles with the markers '[B]' for items being traded or '[S]' for sought after items [44].

*2) Keyword-based post extraction:* In other cases, post contents are sufficient for the researcher's purpose and they will extract posts of interest based on specific keywords. For instance, this is the case in work aimed at understanding hackers' interests and skills in the Internet of Things, based on an analysis of Shodan related posts [5].

*3) Extracting a random sample of posts:* For other types of analysis a random sample is satisfactory, mostly with the aim of constraining the original vast dataset to a smaller set of samples. This is particularly useful for cases where manual annotation is applied, for example prior to training an ML model. Annotation is a slow and expensive process, especially

if multiple annotators are used, so in such cases it may be helpful to reduce the dataset to the smallest usable size, for example using random stratified sampling [62].

*4) Extracting posts or threads based on metadata:* Instead of a random set of samples or posts selected based on specific keywords, some research requires the extraction of a sample based on other criteria, such as posts created by specific authors, within a date range, or threads in a specific bulletin board. For example, Caines et al. extracted posts from selected bulletin boards on HACK FORUMS, such as 'Beginners Hacking' and 'Premium Sellers', filtering for threads with fewer than 20 posts to prepare a training set for ML [10].

*5) Extracting posts based on network properties:* In this case posts are extracted to provide information about the activity of the users who created them. For example, an interaction network based on posting activities can help us to analyse individual users as participants in communities [52].

*6) Extracting specific information:* Other extraction can be analysis-dependent. For example, the approach used for dark jargon discovery required the 10,000 most frequent terms to be extracted from selected dark web forums [60].

### D. Methods used

We studied the tools and methods used to analyse underground and other forums. The latter can be broadly grouped into *qualitative* and *quantitative* methods. We further characterise the quantitative methods according to whether they take a modelling, analytical, or social-network approach. Table I shows these categories along with if the given paper analyses a single forum or multiple forums, and if takes into account the whole forum or a subset.

*1) Qualitative methods:* Several qualitative methods have been used to better understand underground forums. Thematic Analysis involves generating themes that describe the analysed posts [5]; it has been used to discover perceptions of gender [4], to contrast types of online money laundering schemes [43], and for analysing topics discussed by members who cross from one forum to another [21], [48]. Crime script analysis involves breaking complex crime types down into a series of steps. It usually takes a qualitative approach, although some steps may be quantified. An example application is to understand the eWhoring business model [30].

*2) Modelling approaches:* ML can be useful for a variety of automation tasks, particularly data categorisation. The basic workflow involves cleaning and transforming the extracted data to a numerical representation suitable for processing. NLP is used mostly for language analysis tasks, such as classifying posts into pre-defined categories based on linguistic features [10], identifying products [17], discovering topics discussed on forums or by specific users [26], [32], exploring dark jargon [59], [73], and authorship attribution [2].

Supervised classification, which requires a labelled training and test set, has been applied to classify exploits [72], [3], understand and predict private interactions [63], [47], extract cyber threat intelligence [13], [55], [12], [45], [58], [57], identify transactional posts and product types mentioned in

them [53], and discover supply chains as a means of investigating the underground economy[35], [7]. Unsupervised methods may be used to explore forums, for example through clustering [65]. Problems that involve social structure, such as identifying illicit products and key players, may be tackled using classifiers based on graph neural networks [18] and embeddings [76]. Such models may experience concept drift, causing the accuracy of classification results to decline over time [54]. Graph-based ML approaches have also been applied to forum datasets, such as Structured Heterogeneous Information Networks. These can model different object and relationship types for classification tasks, such as for modelling underground forum communities to identify key players [77], or the discovery of cybercrime-suspected threads [38].

*3) Quantitative analytical methods:* Statistical techniques providing valuable insights to underground forums include sequential rule mining to detect trends [40], and group-based trajectory modelling, which groups time-series data to trends over time [27]. These can also be used to examine the signals of trust presented in underground forums and markets [24].

*4) Social Network Analysis:* Underground forums naturally lend themselves to analysis of the communities that form on them. SNA techniques can be used for community detection [22], [52], [41], to understand the macro properties of social networks [44], [52], understand how governance works on forums [46], and identify key players [49], [25]. Further quantitative methods include quantifying activity: extracting prices from forum posts and analysing currency exchanges [53]; counting technical terms and the length of posts on forums [25]; and measuring and visualising forum datasets [15]. Some researchers use mixed methods, incorporating both quantitative and qualitative analysis. Such methods combine both depth and breadth. Examples include applying thematic analysis and NLP techniques to understand perceptions of gender [4], the provisioning of booter services [29], and tracking the evolution of a cybercrime market [70].

### E. Tools and programming languages

We analysed the tools and programming languages reported by researchers for underground forum analysis. Most qualitative analyses use NVivo [30]. ML and NLP tools include NLTK [26], [37], [4], [25], [63], spaCy [26], Tree-Tagger [2], Stanford CoreNLP [75], [17], [77], LibSVM [75], [38], Stanford NER system [17], Keras [72], scikit-learn, Word2Vec [73], and Doc2Vec [18]. Tools used for SNA included NetworkX, while other data analysis tools and programming languages included SQL [5], [4], [70], [14], Python [72], [4], [26], Excel [4], [5], BeautifulSoup HTML parser [72], [63], and Stata [24]. Most of these tools require specialised skills and many require data to be in a specific format. This shows the need for a tool that automates analysis and data exploration, and allows data to be exported in a variety of standard formats.

### III. THE CRIMEBB DATASET

CRIMEBB is a collection of underground forum discussions spanning more than a decade, and a result of the Cambridge

Cybercrime Centre's efforts aimed at enabling cybercrime research at scale [50]. We have been scraping various surface web and dark web forums active since 2002 and make the data available to the research community through data-sharing agreements. At the time of writing, CRIMEBB contains more than 99 million posts from 34 forums; and cybercrime researchers are taking advantage of the dataset: it has been shared with 168 scholars through 48 agreements from 37 institutions in 16 countries (excluding us).

Forums within CRIMEBB are all structured in a similar manner: they can be split to sub-forums or *bulletin boards* that are centred around various technical, non technical, specific hacking related subjects, such as 'Hacking Wireless', and potentially illicit activities, such as 'Make black money', or 'Remote Administration Tools'. *Members* can initiate or contribute to *threads*, which are discussions focusing on a specific subject within the selected bulletin board.

### A. Research ethics

We have approval from the Department's ethics committee for the collection of CRIMEBB and its subsequent sharing with other researchers. The ethics case does not give blanket permission for any and all research on the data. As part of the data sharing agreement, researchers need to ensure that their use of the data does not run counter to local law or the ethics regime at their institution.

### B. A workflow of underground forum analysis

To illustrate the problems we aim to solve with POSTCOG, we demonstrate a typical analysis workflow. The purpose is to understand how a data exploration and analysis toolkit can make the curated crime forum data more accessible to users from different backgrounds, and facilitate collaboration in interdisciplinary teams. The example research scenario is based on collaborative work in which researchers combined criminology and computer science expertise to understand the 'eWhoring' business model [30]. The authors used crime script analysis, extracting 6,519 forum posts written by 2,401 members in 297 threads that provide tutorials relating to eWhoring. The steps of data extraction reported in the original analysis were the following:

- *Data access*: Once researchers have completed the data sharing agreement, they are allowed to download the requested dataset, which they restore to a PostgreSQL database locally on the command line.
- *Data exploration*: Data exploration is then achieved by means of SQL commands. This also involves users familiarising themselves with the database schema.
- *Data selection*: Researchers decide which data to extract, such as posts or threads of interest. Data exploration tasks can be carried out using statistical tools, data visualisation tools, or data science packages from the Python ecosystem, such as *pandas*, which is widely used for data exploration. However users must have a working knowledge of the tools and the data processing pipeline

– that is, passing data to the selected tool, and converting it into the required format for analysis.

In this case study, the researchers extracted threads that contain tutorials from the 'E-Whoring' bulletin board in HACK FORUMS by selecting for the keywords '[TUT]' or 'guide' in the thread titles. They then used a ML classifier to filter out threads that were asking questions, rather than providing tutorials. At the end of these steps, the selected posts were ready and the researchers imported them into NVivo for analysis.

### IV. POSTCOG: A TOOLKIT FOR THE INTERDISCIPLINARY ANALYSIS OF UNDERGROUND FORUMS

In an ideal interdisciplinary setting, researchers from different backgrounds with complementary skills could access and explore curated data together. However, not all researchers will have experience with data science tools or access to computer scientist colleagues. There are also various ways of adding value to forum data, by indexing it for searching, or by integrating NLP classifiers to automatically categorise text, which help researchers filter and extract data in different ways. We developed POSTCOG to abstract the technical details of accessing and exploring curated crime forum data for users, specifically to empower people from non-technical backgrounds.

POSTCOG is a web application currently accessible to users of curated crime forum data. Its development is in its second iteration. The first prototype was built on NodeJS, ExpressJS and PostgreSQL. Based on user testing results (see §V), this technology stack was extended by ReactJS and Elasticsearch in the second iteration to improve performance and the presentation of data filters. The following considerations were taken into account when designing the system:

- *Performance*: One key consideration is the ability to return query results to users in close to real-time.
- *Learning curve reduction*: POSTCOG aims to reduce the time and effort required to become familiar with a dataset, to allow user to instead focus on their research questions.
- *Usability*: The application is straightforward to use and is accessible, as assessed through user testing.
- *Functionality*: Users should be able to filter search results by forum and subforum, date, and NLP tags, and to export results as a CSV file for analysis by other applications (as further detailed in §IV-C).

### A. Data pipeline

Figure 1 shows the main components of POSTCOG and the interaction between POSTCOG and the underlying database. Raw forum data, stored in a PostgreSQL database, is exported to a JSON file. The database can be regularly updated through scraping and have additional labels added by our NLP classifiers. The data are imported into an Elasticsearch cluster using the Logstash data processing pipeline. The presentation of data is handled by the Reactivesearch library that provides UI components to build data driven applications.
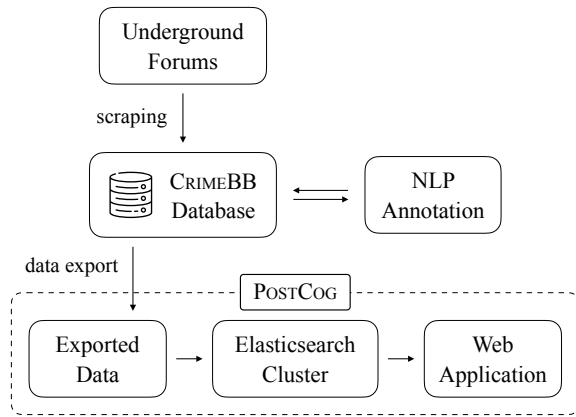
Figure 1: PostCog architecture overview

## B. Index design

A core aspect of PostCog is its data backend, which is implemented using Elasticsearch to allow efficient querying of forum data. In the original data source, forum discussions are stored in a relational database, with post, member, thread and forum information being separated into their dedicated tables [50]. For the purposes of PostCog, an index was designed that contains all of this information. Thus, the export statement responsible for exporting data from the original data source contains data that is the result of joining multiple tables together to ensure the availability of forum, bulletin board, post and NLP label details. The original tables are mapped to a single index and within each document in the index, the properties map to the columns of the original database tables.

## C. Functionality

The primary function of PostCog is to allow users to explore posts within a curated crime forum. Common usage scenarios indicate the most basic form of exploration involves extracting posts discussing a topic of interest.

*1) PostCog dashboard:* Users are required to log in with credentials supplied after signing a data licensing agreement with the Cambridge Cybercrime Centre. The PostCog landing page then provides users a high-level overview of CrimeBB. A dashboard presents descriptive information about the number of posts, number of users, and specifics of forums within the active dataset. This aims to help researchers writing the methods sections of their papers, in which they normally set out a high-level summary of their data sources.

*2) Keyword search:* When a researcher might want to explore IoT hacking over some time period, their first step might be to identify keywords to extract posts containing relevant discussions. PostCog's search functionality supports keyword search by simply typing the search terms. PostCog then displays relevant posts in its results view, along with the titles of the discussion threads, the creation date, forum and bulletin board information and associated NLP labels. Users can search for keywords in both post contents and thread titles. Figure 2 shows the result of a keyword search for the term

'mirai', and is part of a discussion titled '[FREE] World's Largest Net: Mirai Botnet, Client, Echo Loader, CNC source code release.' on Hack Forums within the 'Botnets, IRC Bots, and Zombies' bulletin board. This post announced the release of the source code of Mirai, an Internet of Things botnet subsequently used for many denial-of-service attacks.

*3) Thread view:* Posts are contained in threads, although some threads consist of a single post. Researchers often benefit from being able to view posts as part of a larger context.

*4) Data filters:* Posts can be categorised based on accompanying metadata including the creation time, forum, and bulletin board. Users can constrain search results based on these attributes. For example, users might be interested in a specific bulletin board from one forum, rather than searching all of the 34 forums in CrimeBB. Additionally, users might want to focus on a specific date range, such as the period immediately following a high profile event.

*5) NLP labels:* PostCog allows modular inclusion of NLP classifiers. Users can filter Hack Forums posts using classifiers built by Caines et al. [10] which annotate the nature and intent of a post, and a crime type classifier built by Atondo Siu et al. [61]. NLP labels are currently available for Hack Forums only, although we are evaluating whether the classifiers generalise to other forums. The first classifier categorises posts based on their 'type' (whether it is a comment, a tutorial, etc.) and their 'intent' (whether it was a neutral post, or features aggressive language). The crime type classifier tags posts discussing unauthorised access to systems, bots and malware, currency exchange, denial of service attacks, identity theft, spam, trading credentials, and VPN and hosting services. Posts not overtly discussing illicit activities are tagged 'not criminal'. PostCog integrates these labels so users can select specific categories. The example in Figure 2 has been labelled as a comment. It has also been mislabelled as aggressive. As detailed below, users can flag mislabelled posts to improve the performance of our classifiers. As the NLP label functionality is modular, we can integrate additional classifiers.

*6) Offline analysis and data export:* Depending on the type of analysis, users might wish to export the filtered search results for analysis with a specialised tool, so export to a CSV file is supported.

*7) Feedback and reporting:* Due to inherent characteristics of ML models and domain complexity, not all NLP labels will reflect the true nature of posts. To improve the models, annotations and user experience, PostCog allows users to report incorrect labels. They can also flag posts containing hate speech or hateful content: to help users, we provide the United Nations' hate speech definition [69]. Annotation of posts – whether author intent, crime type or hate speech – is a nontrivial task. ML models need constant human feedback for training, just as with spam filters.

*8) Information source:* PostCog also integrates information about the underlying data. Intended for more advanced users, this includes information about the database schema, how users can set it up locally if they wish to, and information about underground forums.

| Date: | 2016-09-30 | Bulletin board: | Botnets, IRC Bots, and Zombies | Forum: | Hack Forums |
|---|---|---|---|---|---|
| Intent: | aggression | Post type: | comment | | |

Preface Greetz everybody. When I first go in DDoS industry, I wasn't planning on staying in it long. I made my money, there's lots of eyes looking at IOT now, so it's time to GTFO. However, I know every skid and their mama, it's their wet dream to have something besides qbot. So today, I have an amazing release for you. With Mirai, I usually pull max 380k bots from telnet alone. However, after the Kreb DDoS, ISPs been slowly shutting down and cleaning up their act. Today, max pull is about 300K bots, and dropping. So, I am your senpai, and I will treat you real nice, my hf-chan. And to everyone that thought they were doing anything by hitting my CNC, I had good laughs, this bot uses domain for CNC. It takes 60 seconds for all bots to reconnect, lol Also, shoutout to this blog post by malwaremustdie...

| View post in thread | Flag incorrect post class | Flag hate speech |
|---|---|---|

Figure 2: An example post on HACK FORUMS announcing the release of Mirai Botnet's source code, displayed by POSTCOG.

### D. Example analysis workflow revisited

We return to the example of extracting eWhoring tutorials introduced in §III-B to highlight the differences between the original approach and the POSTCOG route. The steps of extracting the relevant data using POSTCOG are as follows:

- Select 'E-Whoring' board and the 'HACK FORUMS' forum.
- Enter keyword '[TUT]' and 'guide' in the 'Search for keyword in thread titles' search bar. Users can also utilise the NLP label filters by selecting the 'Tutorial' option in the 'Post type' filter. This displays posts classified as providing tutorials, rather than asking questions.

The eWhoring tutorials will be displayed, and can be downloaded to a CSV file for further work. For the example provided, the data would be imported into NVivo for qualitative crime script analysis.

### V. PROTOTYPE EVALUATION

The prototype was evaluated by two rounds of user testing. As no existing tools met the need that POSTCOG was designed to meet, we could not do a comparative analysis. Instead, for both the initial and subsequent evaluations, we used the standard usability research method 'Thinking aloud testing' [20]. Representative users were asked to use the system and describe their thoughts as they carried out tasks. These included specific and open-ended tasks to evaluate how well the prototype satisfies standard usability heuristics for user interface design, and as a result how usable the system is for users in different disciplines. Direct observation and dialogue with test subjects provided quick qualitative feedback on both the prototype's positive aspects (required information is easy to find, feature is useful, etc.) and its negative ones (navigation issues, insufficient information, high cognitive load, etc.), which informed further design and development.

### A. Research ethics

Approval was obtained from the department's ethics committee for the usability studies. Participants were not offered a monetary incentive, but were offered a copy of the published research when available. The provision of the tool is also a tangible benefit for participants. Participants were also provided with the opportunity to be listed in the acknowledgement section, should they wish to receive recognition.

### B. Process and participants

A total of two user testing sessions took place. The first session was aimed at assessing the usability of the prototype. The second round tested the second iteration of the prototype with improvements and changes based on feedback from the first session. Due to the pandemic, the usability tests were conducted online using Zoom. Participants accessed POSTCOG using a browser on their own computer through a screen-sharing session. The calls were not recorded, but the researcher took detailed notes. Participants spent around one hour with the researcher during each test.

In total, four users participated in both rounds of user testing. All those that were invited to participate in the evaluation agreed to take part. The participants were recruited from among the active users of the CRIMEBB dataset, based on their research background and representativeness of POSTCOG's user base; the evaluation involved users from both social science and computer science.

The test subjects possessed varying levels of technical skill and experience with forum analysis. One had prior analysis experience with the specific dataset and with various data extraction and analysis tools, such as SQL and Python. The other three users had experience in applying qualitative methods to data extracted from CRIMEBB.

### C. First evaluation

The first evaluation focuses on assessing whether and how the prototype fulfils its high level goals. This subsection provides a description of our research questions and associated tasks users were asked to perform using the prototype.

*1) What are users' perceptions of the prototype?:* The task associated with this question was an open-ended one to assess whether users find POSTCOG intuitive, which allowed us to

gather useful feedback about various facets of the web application including the information architecture and navigation, user interface elements, the presentation of search results and data filters, and the performance aspects of querying large datasets. Users were asked to explore POSTCOG including its main page that contains a dashboard for a quick overview of the data, its data explorer, and its 'about' page that describes the dataset and NLP labels in greater detail.

→ *Result:* All users agreed that the dashboard on the main page was a useful feature that provides an overview of the dataset and that overall the information presented on this page is structured clearly. However the users pointed out that the 'about' page lacked clear navigation and would require content changes. One of the users also struggled to create filters to select forums and bulletin boards, and pointed out that examples of how to use the data explorer would be useful. We realised that the filters were not presented in an intuitive manner and that suggested substantial changes were required. All users noted the filters representing NLP labels were not clearly presented, and more information was required to introduce them. Finally, while the prototype aimed to provide a minimal working version of POSTCOG, the tests highlighted the need for scalability of the data explorer. This led us to seek out technologies that support near real-time large scale data analysis for our second iteration.

*2) Can users explore forum posts and filter for posts with ease?:* The users were asked to carry out tasks that reflect real-world use to discover if the filters were intuitive, to assess how effectively they could perform specific exploration tasks such as searching for a particular keyword on a specific forum in the defined date range. For example, "Select a sample of posts from HACK FORUMS. How many posts are returned? Order the posts in descending date order. What date is the first post from? What board is this post associated with?", or "Search for posts that contain the term 'malware'. How many posts are returned? Name some threads these posts are associated with." The presentation of NLP labels was evaluated through tasks such as "Identify the main 'post type' categories, and summarise in your own words what they mean". Complementing the open-ended UI exploration task, this part of the evaluation identified issues with specific UI elements.

→ *Result:* All users reported the date picker was cumbersome to use, and example queries and explanations of filters should be incorporated in the form of additional UI elements. Users also suggested useful additional functionality for data exploration, such as a query builder feature, which can be incorporated in subsequent phases of the development.

*3) Can users reconstruct conversation threads with ease?:* Such tasks included "View the full thread of posts for the first post returned from Task n. Order the posts in the thread by timestamp. Identify the first message in the thread. What is the time difference between the first post and the next?"

→ *Result:* All users reported that being able to see full conversation threads rather than just the individual posts was a highly beneficial feature, as analysis is supported by access to contextual information.

*4) Can users save posts as CSV files for further analysis?:* To test this, users were asked to download selected posts to a CSV file, open it, and examine the contents including the header, to ensure information related to posts is presented in such a way as to enable further analysis.

→ *Result:* All users said that this feature was helpful and did not report any issues with invoking it.

*5) Can users easily provide feedback on identifying posts that contain hate speech or that have been labelled incorrectly?:* "Select a sample of posts that are tagged as 'Tutorials'. How many posts are returned? Read a few posts. Are there any posts where you disagree with the label, i.e. they are incorrectly tagged as 'Tutorial'? If so, flag the post."

→ *Result:* This exercise highlighted another area where the prototype needed major revisions. The users could not locate relevant filters easily and found that the information presented to them was not sufficient to make a decision about labels.

*6) Can users easily access the dashboard, which provides an overview of the data in the currently active dataset? As part of this, can they find information related to the scale of the data, and relevant metrics such as forum size?:* The users were asked to perform tasks such as "Identify English language forums", "Which forum has the highest number of posts?", "Observe how the number of posts over time has changed".

→ *Result:* While all users successfully completed their tasks, they suggested improvements. These involved providing more information on forums, for example whether they are dark web or surface web forums, or turning the names of forums into hyperlinks pointing to more detailed explanations.

*D. Second evaluation*

Overall, results from the first round suggest that while users found the prototype useful, it required major improvement. The prototype was duly revised and a second prototype was built. In the second round of user testing we turned to specific underground forum analysis scenarios found in the literature to explore whether users can perform these analysis tasks with ease with POSTCOG. In particular, users were asked to repeat the analysis workflow introduced in Section III-B to assess whether POSTCOG indeed reduces the steep learning curve in a qualitative analysis task. Users were asked to extract eWhoring tutorials from the HACK FORUMS dataset's 'E-Whoring' bulletin board.

→ *Result:* The second prototype was tested with the same users, all of whom completed the practical task and successfully extracted eWhoring tutorials from HACK FORUMS. They used right filters effectively and reported that the interface was intuitive. The exercise emphasised the need for advanced search functionality and the ability to search for multiple keywords. Users also wanted to combine multiple keywords using logical operators. Further requirements that emerged and that have now been implemented include:

- Navigation: we added a dedicated page to present information related to the NLP classifiers in POSTCOG.
- UI elements: we fixed date picker usability issues, a missing contact form, and missing information on UI

elements including table legends or pop-up dialogs which explain what dark web forums are in the overview table.

- New functionality: we added links to publications providing further information on forum analysis, associating the links the relevant forum in the forums overview table.

## VI. RELATED WORK

We outline related work regarding underground forums in §II. To the best of our knowledge, POSTCOG is the first tool designed to enable researchers understand the cybercrime ecosystem. We were unable to find similar tools designed to enable research at scale into forum communities more generally. The CRIMEBB dataset has many unique properties in the way it is structured and the nature of its content, and we used our domain expertise and knowledge of researcher interests to develop a bespoke user interface to a search engine built from open-source components.

There is a body of work relating to the usability of search engines [16], including for young users [23] and blind users [9]. We searched more specifically for similar tools designed and evaluated for research communities. We only found tools designed for medical research, such as NLP tools for extracting information from clinical notes [28], [67], [66] and electronic health records [33]. Large biomedical projects sometimes offer bespoke search engines for domain-specific literature [42], [36], [56] as well as data output from bioinformatics research and text mining [11], [6], [39].

NLP methods have been used extensively to help end users access relevant information more quickly online: indeed, research into information retrieval underpins and predates search engines such as Google, Bing and Yahoo. Many other types of website include search facilities, from online retailers to streaming services to social networks, and there are search engines serving particular domains. We are not aware of existing search engines for cybercrime forum data or for underground forum data more generally, and therefore believe that POSTCOG is the first of its kind.

## VII. DISCUSSION

Analysing discussions and users on underground forums provides valuable insights into cybercrime, underground economies, and the social networks and communities formed on these forums. Such analysis is inherently interdisciplinary, and requires better ways for researchers from non-technical backgrounds to undertake research based on their domain knowledge of the area but without having to acquire deep knowledge of complex data-science tools.

Our first usability studies demonstrated challenges with our initial POSTCOG prototype, leading to a much improved product. As ever with web development, we will continue to work on new features, along with new requirements identified by future user testing and feedback. We will continue to incorporate new data, and new annotation methods based on classifiers relating to new and existing data. Readers may access our existing curated forum data using the current version of POSTCOG by completing a data licensing agreement.

POSTCOG is still in its early stages. We will continue to develop and evaluate it with more participants, as well as expanding it to index content in EXTREMEBB – a newly established dataset covering extremist forums [71]. The continued collection of both cybercrime and extremist forum data requires constant maintenance to defeat anti-scraping countermeasures implemented by administrators [68]. Besides forums, we also plan to expand POSTCOG to work with other cybercrime and extremist material including Telegram and Discord channels, which we also collect, and to anonymous underground markets too.

We anticipate a flagging system for incorrect labels will help improve our classifiers' performance. We aim to adapt the NLP classifiers to other languages and forums, including those used by violent online political extremists. In the long term, we also plan to support more analysis methods such as extracting random samples satisfying particular constraints, and interactive interfaces for social network analysis, as well as indicators of economic activity. POSTCOG's modular design supports further evolution, while the use of Elasticsearch helps ensure it is scalable as our datasets expand. Our long-term aim is to develop POSTCOG into a research engine for all kinds of online wickedness.

As cybercrime is adversarial, forum users may try to impede our efforts. There are already poisoning attacks against NLP models for which we need to implement defences; these include sanitising data to minimise the likelihood that users can use text-encoding attacks to disrupt our ML models [8].

## VIII. CONCLUSION

To develop better tools for underground forum analysis, we surveyed relevant research then studied steps involved in the analysis, data volume, and tools used for each step. From this we developed a web application, POSTCOG, which functions as a search engine for underground forums. It provides users an intuitive way to explore, filter and export forum posts that have been scraped and curated for offline analysis. We evaluated POSTCOG with the help of users from an interdisciplinary cybercrime research group, enabling us to extend and improve it. Its use with the CRIMEBB dataset is subject to a license available at zero cost to bona fide academic researchers.

## IX. ACKNOWLEDGEMENTS

REFERENCES

[1] S. Afroz, V. Garg, D. McCoy, and R. Greenstadt, "Honor Among Thieves: A Common's Analysis of Cybercrime Economies," in *Proc. IEEE eCrime Researchers Summit (eCrime)*, 2013, pp. 1–11.

[2] S. Afroz, A. C. Islam, A. Stolerman, R. Greenstadt, and D. McCoy, "Doppelgänger Finder: Taking Stylometry to the Underground," in *Proc. IEEE Symposium on Security and Privacy (S&P)*, 2014, pp. 212–226.

[3] B. Ampel, S. Samtani, H. Zhu, S. Ullman, and H. Chen, "Labeling Hacker Exploits for Proactive Cyber Threat Intelligence: A Deep Transfer Learning Approach," in *Proc. IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2020, pp. 1–6.

[4] M. Bada, Y. T. Chua, B. Collier, and I. Pete, "Exploring Masculinities and Perceptions of Gender in Online Cybercrime Subcultures," in *Cybercrime in Context: The Human Factor in Victimization, Offending, and Policing*, 2021, pp. 237–257.

[5] M. Bada and I. Pete, "An Exploration of the Cybercrime Ecosystem Around Shodan," in *Proc. IEEE International Conference on Internet of Things: Systems, Management and Security (IOTSMS)*, 2020, pp. 1–8.

[6] A. P. Bento, A. Gaulton, A. Hersey, L. J. Bellis, J. Chambers, M. Davies, F. A. Krüger, Y. Light, L. Mak, S. McGlinchey *et al.*, "The ChEMBL Bioactivity Database: An Update," *Nucleic Acids Research*, vol. 42, no. D1, pp. D1083–D1090, 2014.

[7] R. Bhalerao, M. Aliapoulios, I. Shumailov, S. Afroz, and D. McCoy, "Mapping the Underground: Supervised Discovery of Cybercrime Supply Chains," in *Proc. IEEE Symposium on Electronic Crime Research (eCrime)*, 2019, pp. 1–16.

[8] N. Boucher, I. Shumailov, R. Anderson, and N. Papernot, "Bad Characters: Imperceptible NLP Attacks," *arXiv preprint arXiv:2106.09898*, 2021.

[9] M. Buzzi, P. Andronico, and B. Leporini, "Accessibility and Usability of Search Engine Interfaces: Preliminary Testing," in *Proc. ERCIM UI4ALL Workshop*, 2004.

[10] A. Caines, S. Pastrana, A. Hutchings, and P. J. Buttery, "Automatically Identifying the Function and Intent of Posts in Underground Forums," *Crime Science*, vol. 7, no. 1, pp. 1–14, 2018.

[11] F. Cunningham, M. R. Amode, D. Barrell, K. Beal, K. Billis, S. Brent, D. Carvalho-Silva, P. Clapham, G. Coates, S. Fitzgerald *et al.*, "Ensembl 2015," *Nucleic Acids Research*, vol. 43, no. D1, pp. D662–D669, 2015.

[12] A. Deb, K. Lerman, and E. Ferrara, "Predicting Cyber-Events by Leveraging Hacker Sentiment," *Information*, vol. 9, no. 11, p. 280, 2018.

[13] I. Deliu, C. Leichter, and K. Franke, "Extracting Cyber Threat Intelligence From Hacker Forums: Support Vector Machines versus Convolutional Neural Networks," in *Proc. IEEE International Conference on Big Data (BigData)*, 2017, pp. 3648–3656.

[14] ——, "Collecting Cyber Threat Intelligence from Hacker Forums via a Two-Stage, Hybrid Process Using Support Vector Machines and Latent Dirichlet Allocation," in *Proc. IEEE International Conference on Big Data (BigData)*, 2018, pp. 5008–5013.

[15] P.-Y. Du, N. Zhang, M. Ebrahimi, S. Samtani, B. Lazarine, N. Arnold, R. Dunn, S. Suntwal, G. Angeles, R. Schweitzer *et al.*, "Identifying, Collecting, and Presenting Hacker Community Data: Forums, IRC, Carding Shops, and DNMs," in *Proc. IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2018, pp. 70–75.

[16] D. Dudek, A. Mastora, and M. Landoni, "Is Google the answer? A Study into Usability of Search Engines," *Library Review*, vol. 56, no. 3, pp. 224–233, 2007.

[17] G. Durrett, J. K. Kummerfeld, T. Berg-Kirkpatrick, R. S. Portnoff, S. Afroz, D. McCoy, K. Levchenko, and V. Paxson, "Identifying Products in Online Cybercrime Marketplaces: A Dataset for Fine-grained Domain Adaptation," *arXiv preprint arXiv:1708.09609*, 2017.

[18] Y. Fan, Y. Ye, Q. Peng, J. Zhang, Y. Zhang, X. Xiao, C. Shi, Q. Xiong, F. Shao, and L. Zhao, "Metagraph Aggregated Heterogeneous Graph Neural Network for Illicit Traded Product Identification in Underground Market," in *Proc. IEEE International Conference on Data Mining (ICDM)*, 2020, pp. 132–141.

[19] Y. Fang, Y. Guo, C. Huang, and L. Liu, "Analyzing and Identifying Data Breaches in Underground Forums," *IEEE Access*, vol. 7, pp. 48 770–48 777, 2019.

[20] M. E. Fonteyn, B. Kuipers, and S. J. Grobe, "A Description of Think Aloud Method and Protocol Analysis," *Qualitative Health Research*, vol. 3, no. 4, pp. 430–441, 1993.

[21] R. Frank, M. Thomson, A. Mikhaylov, and A. J. Park, "Putting All Eggs in a Single Basket: A Cross-community Analysis of 12 Hacking Forums," in *Proc. IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2018, pp. 136–141.

[22] V. Garg, S. Afroz, R. Overdorf, and R. Greenstadt, "Computer-Supported Cooperative Crime," in *Proc. International Conference on Financial Cryptography and Data Security (FC)*, 2015, pp. 32–43.

[23] T. Gossen, J. Hempel, and A. Nürnberger, "Find it if You Can: Usability Case Study of Search Engines for Young Users," *Personal and Ubiquitous Computing*, vol. 17, no. 8, pp. 1593–1603, 2013.

[24] T. J. Holt, O. Smirnova, and A. Hutchings, "Examining Signals of Trust in Criminal Markets Online," *Journal of Cybersecurity*, vol. 2, no. 2, pp. 137–145, 2016.

[25] C. Huang, Y. Guo, W. Guo, and Y. Li, "HackerRank: Identifying Key Hackers in Underground Forums," *International Journal of Distributed Sensor Networks*, vol. 17, no. 5, pp. 1–12, 2021.

[26] J. Hughes, S. Aycock, A. Caines, P. Buttery, and A. Hutchings, "Detecting Trending Terms in Cybersecurity Forum Discussions," in *Proc. Workshop on Noisy User-generated Text (WNUT)*, 2020, pp. 107–115.

[27] J. Hughes, B. Collier, and A. Hutchings, "From Playing Games to Committing Crimes: A Multi-technique Approach to Predicting Key Actors on an Online Gaming Forum," in *Proc. IEEE Symposium on Electronic Crime Research (eCrime)*, 2019, pp. 1–12.

[28] G. Hultman, R. McEwan, S. Pakhomov, E. Lindemann, S. Skube, and G. B. Melton, "Usability Evaluation of an Unstructured Clinical Document Query Tool for Researchers," *AMIA Summits on Translational Science Proceedings*, vol. 2018, p. 84, 2018.

[29] A. Hutchings and R. Clayton, "Exploring the Provision of Online Booter Services," *Deviant Behavior*, vol. 37, no. 10, pp. 1163–1178, 2016.

[30] A. Hutchings and S. Pastrana, "Understanding eWhoring," in *Proc. IEEE European Symposium on Security and Privacy (EuroS&P)*, 2019, pp. 201–214.

[31] J. W. Johnsen and K. Franke, "Identifying Central Individuals in Organised Criminal Groups and Underground Marketplaces," in *Proc. International Conference on Computational Science (ICCS)*, 2018, pp. 379–386.

[32] ——, "Identifying Proficient Cybercriminals Through Text and Network Analysis," in *Proc. IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2020, pp. 1–7.

[33] D. R. Kaufman, B. Sheehan, P. Stetson, A. R. Bhatt, A. I. Field, C. Patel, and J. M. Maisel, "Natural Language Processing–Enabled and Conventional Data Capture Methods for Input to Electronic Health Records: a Comparative Usability Study," *JMIR Medical Informatics*, vol. 4, no. 4, p. e5544, 2016.

[34] B. Krebs. (2019) No Jail Time for "WannaCry Hero". https://krebsonsecurity.com/2019/07/no-jail-time-for-wannacry-hero/. Online; accessed 20 April 2022.

[35] M. K. Kumar and D. K. Bhargavi, "An Effective Study on Data Science Approach to Cybercrime Underground Economy Data," *Journal of Engineering, Computing and Architecture*, vol. 10, no. 1, pp. 148–158, 2020.

[36] M. Levchenko, Y. Gou, F. Graef, A. Hamelers, Z. Huang, M. Ide-Smith, A. Iyer, O. Kilian, J. Katuri, J.-H. Kim *et al.*, "Europe PMC in 2017," *Nucleic Acids Research*, vol. 46, no. D1, pp. D1254–D1260, 2018.

[37] Y. Li, J. Cheng, C. Huang, Z. Chen, and W. Niu, "NEDetector: Automatically Extracting Cybersecurity Neologisms from Hacker Forums," *Journal of Information Security and Applications*, vol. 58, p. 102784, 2021.

[38] J. Liu, "Automatic Detection of Cybercrime-suspected Threads in Online Underground Forums," *Graduate Theses, Dissertations, and Problem Reports. 7204*, 2018.

[39] Y. Liu, Y. Liang, and D. Wishart, "PolySearch2: A Significantly Improved Text-mining System for Discovering Associations between Human Diseases, Genes, Drugs, Metabolites, Toxins and More," *Nucleic Acids Research*, vol. 43, no. W1, pp. W535–W542, 2015.

[40] E. Marin, M. Almukaynizi, E. Nunes, J. Shakarian, and P. Shakarian, "Predicting Hacker Adoption on Darkweb Forums Using Sequential Rule Mining," in *Proc. IEEE International Conference on Parallel Distributed Processing with Applications (ISPDC)*, 2018, pp. 1183–1190.

[41] E. Marin, M. Almukaynizi, E. Nunes, and P. Shakarian, "Community Finding of Malware and Exploit Vendors on Darkweb Marketplaces," in *Proc. IEEE International Conference on Data Intelligence and Security (ICDIS)*, 2018, pp. 81–84.

[42] J. R. McEntyre, S. Ananiadou, S. Andrews, W. J. Black, R. Boulder-stone, P. Buttery, D. Chaplin, S. Chevuru, N. Cobley, L.-A. Coleman *et al.*, "UKPMC: A Full Text Article Resource for the Life Sciences," *Nucleic Acids Research*, vol. 39, pp. D58–D65, 2010.

[43] A. Mikhaylov and R. Frank, "Cards, Money and Two Hacking Forums: An Analysis of Online Money Laundering Schemes," in *Proc. IEEE European Intelligence and Security Informatics Conference (EISIC)*, 2016, pp. 80–83.

[44] M. Motoyama, D. McCoy, K. Levchenko, S. Savage, and G. M. Voelker, "An Analysis of Underground Forums," in *Proc. ACM Internet Measurement Conference (IMC)*, 2011, pp. 71–80.

[45] E. Nunes, A. Diab, A. Gunn, E. Marin, V. Mishra, V. Paliath, J. Robertson, J. Shakarian, A. Thart, and P. Shakarian, "Darknet and Deepnet Mining for Proactive Cybersecurity Threat Intelligence," in *Proc. IEEE Conference on Intelligence and Security Informatics (ISI)*, 2016, pp. 7–12.

[46] M. Odabaş, T. J. Holt, and R. L. Breiger, "Markets as Governance Environments for Organizations at the Edge of Illegality: Insights From Social Network Analysis," *American Behavioral Scientist*, vol. 61, no. 11, pp. 1267–1288, 2017.

[47] R. Overdorf, C. Troncoso, R. Greenstadt, and D. McCoy, "Under the Underground: Predicting Private Interactions in Underground Forums," *arXiv preprint arXiv:1805.04494*, 2018.

[48] A. J. Park, R. Frank, A. Mikhaylov, and M. Thomson, "Hackers Hedging Bets: A Cross-Community Analysis of Three Online Hacking Forums," in *Proc. IEEE International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 2018, pp. 798–805.

[49] S. Pastrana, A. Hutchings, A. Caines, and P. Buttery, "Characterizing Eve: Analysing Cybercrime Actors in a Large Underground Forum," in *Proc. International Symposium on Research in Attacks, Intrusions, and Defenses (RAID)*, 2018, pp. 207–227.

[50] S. Pastrana, D. R. Thomas, A. Hutchings, and R. Clayton, "CrimeBB: Enabling Cybercrime Research on Underground Forums at Scale," in *Proc. World Wide Web Conference (WWW)*, 2018, pp. 1845–1854.

[51] I. Pete and Y. T. Chua, "An Assessment of the Usability of Cybercrime Datasets," in *Proc. USENIX Workshop on Cyber Security Experimentation and Test (USENIX CSET)*, 2019.

[52] I. Pete, J. Hughes, Y. T. Chua, and M. Bada, "A Social Network Analysis and Comparison of Six Dark Web Forums," in *Proc. IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, 2020, pp. 484–493.

[53] R. S. Portnoff, S. Afroz, G. Durrett, J. K. Kummerfeld, T. Berg-Kirkpatrick, D. McCoy, K. Levchenko, and V. Paxson, "Tools for Automated Analysis of Cybercriminal Markets," in *Proc. World Wide Web Conference (WWW)*, 2017, pp. 657–666.

[54] A. L. Queiroz, B. Keegan, and S. Mckeever, "Moving Targets: Addressing Concept Drift in Supervised Models for Hacker Communication Detection," in *Proc. IEEE International Conference on Cyber Security and Protection of Digital Services (Cyber Security)*, 2020, pp. 1–7.

[55] A. L. Queiroz, S. Mckeever, and B. Keegan, "Detecting Hacker Threats: Performance of Word and Sentence Embedding Models in Identifying Hacker Communications," in *Proc. International Conference on Artificial Intelligence and Computer Science (AICS)*, 2019, pp. 116–127.

[56] D. Rebholz-Schuhmann, J.-H. Kim, Y. Yan, A. Dixit, C. Friteyre, R. Hoehndorf, R. Backofen, and I. Lewin, "Evaluation and Cross-Comparison of Lexical Entities of Biological Interest (LexEBI)," *PLoS One*, vol. 8, no. 10, pp. 1–15, 2013.

[57] A. Sapienza, A. Bessi, S. Damodaran, P. Shakarian, K. Lerman, and E. Ferrara, "Early Warnings of Cyber Threats in Online Discussions," in *Proc. IEEE International Conference on Data Mining Workshops (ICDMW)*, 2017, pp. 667–674.

[58] A. Sapienza, S. K. Ernala, A. Bessi, K. Lerman, and E. Ferrara, "DISCOVER: Mining Online Chatter for Emerging Cyber Threats," in *Proc. World Wide Web Conference (WWW)*, 2018, pp. 983–990.

[59] D. Seyler, W. Liu, X. Wang, and C. Zhai, "Towards Dark Jargon Interpretation in Underground Forums," in *Proc. European Conference on Information Retrieval (ECIR)*, 2021, pp. 393–400.

[60] D. Seyler, W. Liu, Y. Zhang, X. Wang, and C. Zhai, "DarkJargon.Net: A Platform for Understanding Underground Conversation with Latent Meaning," in *Proc. International Conference on Research and Development in Information Retrieval (SIGIR)*, 2021, pp. 2526–2530.

[61] G. A. Siu, B. Collier, and A. Hutchings, "Follow the Money: The Relationship between Currency Exchange and Illicit Behaviour in an Underground Forum," in *Proc. IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, 2021, pp. 191–201.

[62] Z. Sun, A. Oest, P. Zhang, C. Rubio-Medrano, T. Bao, R. Wang, Z. Zhao, Y. Shoshitaishvili, A. Doupé, G.-J. Ahn *et al.*, "Having Your Cake and Eating It: An Analysis of Concession-Abuse-as-a-Service," in *Proc. USENIX Security Symposium (USENIX Security)*, 2021, pp. 4169–4186.

[63] Z. Sun, C. E. Rubio-Medrano, Z. Zhao, T. Bao, A. Doupé, and G.-J. Ahn, "Understanding and Predicting Private Interactions in Underground Forums," in *Proc. Conference on Data and Application Security and Privacy (CODASPY)*, 2019, pp. 303–314.

[64] S. Sundaresan, D. McCoy, S. Afroz, and V. Paxson, "Profiling Underground Merchants Based on Network Behavior," in *Proc. IEEE Symposium on Electronic Crime Research (eCrime)*, 2016, pp. 1–9.

[65] J. Tachaiya, J. Gharibshah, E. E. Papalexakis, and M. Faloutsos, "RThread: A Thread-centric Analysis of Security Forums," in *Proc. IEEE International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 2020, pp. 1–5.

[66] G. Trivedi, E. R. Dadashzadeh, R. M. Handzel, W. W. Chapman, S. Visweswaran, and H. Hochheiser, "Interactive NLP in Clinical Care: Identifying Incidental Findings in Radiology Reports," *Applied Clinical Informatics*, vol. 10, no. 04, pp. 655–669, 2019.

[67] G. Trivedi, P. Pham, W. W. Chapman, R. Hwa, J. Wiebe, and H. Hochheiser, "NLPReViz: an Interactive Tool for Natural Language Processing on Clinical Text," *Journal of the American Medical Informatics Association*, vol. 25, no. 1, pp. 81–87, 2018.

[68] K. Turk, S. Pastrana, and B. Collier, "A Tight Scrape: Methodological Approaches to Cybercrime Research Data Collection in Adversarial Environments," in *Proc. IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, 2020, pp. 428–437.

[69] United Nations Strategy and Plan of Action on Hate Speech. https://www.un.org/en/genocideprevention/hate-speech-strategy.shtml. Online; accessed 20 April 2022.

[70] A. V. Vu, J. Hughes, I. Pete, B. Collier, Y. T. Chua, I. Shumailov, and A. Hutchings, "Turning Up the Dial: The Evolution of a Cybercrime Market Through Set-up, Stable, and Covid-19 Eras," in *Proc. ACM Internet Measurement Conference (IMC)*, 2020, pp. 551–566.

[71] A. V. Vu, L. Wilson, Y. T. Chua, I. Shumailov, and R. Anderson, "ExtremeBB: Enabling Large-Scale Research into Extremism, the Manosphere and Their Correlation by Online Forum Data," *arXiv preprint arXiv:2111.04479*, 2021.

[72] R. Williams, S. Samtani, M. Patton, and H. Chen, "Incremental Hacker Forum Exploit Collection and Classification for Proactive Cyber Threat Intelligence: An Exploratory Study," in *Proc. IEEE International Conference on Intelligence and Security Informatics (ISI)*, 2018, pp. 94–99.

[73] K. Yuan, H. Lu, X. Liao, and X. Wang, "Reading Thieves' Cant: Automatically Identifying and Understanding Dark Jargons from Cybercrime Marketplaces," in *Proc. USENIX Security Symposium (USENIX Security)*, 2018, pp. 1027–1041.

[74] W. T. Yue, Q.-H. Wang, and K.-L. Hui, "See No Evil, Hear No Evil? Dissecting the Impact of Online Hacker Forums," *Mis Quarterly*, vol. 43, no. 1, p. 73, 2019.

[75] Y. Zhang, Y. Fan, S. Hou, J. Liu, Y. Ye, and T. Bourlai, "iDetector: Automate Underground Forum Analysis Based on Heterogeneous Information Network," in *Proc. IEEE International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 2018, pp. 1071–1078.

[76] Y. Zhang, Y. Fan, Y. Ye, L. Zhao, and C. Shi, "Key Player Identification in Underground Forums over Attributed Heterogeneous Information Network Embedding Framework," in *Proc. International Conference on Information and Knowledge Management (CIKM)*, 2019, pp. 549–558.

[77] Y. Zhang, Y. Fan, Y. Ye, L. Zhao, J. Wang, Q. Xiong, and F. Shao, "KADetector: Automatic Identification of Key Actors in Online Hack Forums Based on Structured Heterogeneous Information Network," in *Proc. IEEE International Conference on Big Knowledge (ICBK)*, 2018, pp. 154–161.