



Contents lists available at ScienceDirect

Digital Investigation

journal homepage: www.elsevier.com/locate/diin

Criminal motivation on the dark web: A categorisation model for law enforcement

Janis Dalins^{a, b, *}, Campbell Wilson^a, Mark Carman^a

^a Faculty of Information Technology, Monash University, 900 Princes Highway, Caulfield East, Victoria, Australia

^b Australian Federal Police, 383 La Trobe Street, Melbourne, Victoria, Australia

ARTICLE INFO

Article history:

Received 12 October 2017

Received in revised form

4 December 2017

Accepted 5 December 2017

Available online xxx

Keywords:

Dark web

Computer forensics

Conceptual models

Focused crawls

Machine learning

Child pornography

Tor motivation model

ABSTRACT

Research into the nature and structure of 'Dark Webs' such as Tor has largely focused upon manually labelling a series of crawled sites against a series of categories, sometimes using these labels as a training corpus for subsequent automated crawls. Such an approach is adequate for establishing broad taxonomies, but is of limited value for specialised tasks within the field of law enforcement. Contrastingly, existing research into illicit behaviour online has tended to focus upon particular crime types such as terrorism. A gap exists between taxonomies capable of holistic representation and those capable of detailing criminal behaviour. The absence of such a taxonomy limits interoperability between agencies, curtailing development of standardised classification tools.

We introduce the Tor-use Motivation Model (TMM), a two-dimensional classification methodology specifically designed for use within a law enforcement context. The TMM achieves greater levels of granularity by explicitly distinguishing site content from motivation, providing a richer labelling schema without introducing inefficient complexity or reliance upon overly broad categories of relevance. We demonstrate this flexibility and robustness through direct examples, showing the TMM's ability to distinguish a range of *unethical* and *illegal* behaviour without bloating the model with unnecessary detail.

The authors of this paper received permission from the Australian government to conduct an unrestricted crawl of Tor for research purposes, including the gathering and analysis of illegal materials such as child pornography. The crawl gathered 232,792 pages from 7651 Tor virtual domains, resulting in the collation of a wide spectrum of materials, from illicit to downright banal. Existing conceptual models and their labelling schemas were tested against a small sample of gathered data, and were observed to be either overly prescriptive or vague for law enforcement purposes - particularly when used for prioritising sites of interest for further investigation.

In this paper we deploy the TMM by manually labelling a corpus of over 4000 unique Tor pages. We found a network impacted (but not dominated) by illicit commerce and money laundering, but almost completely devoid of violence and extremism. In short, criminality on this 'dark web' is based more upon greed and desire, rather than any particular political motivations.

© 2018 Elsevier Ltd. All rights reserved.

Introduction

Exhaustive examinations of networks such as the worldwide web (WWW) have been shown to be expensive, inefficient and ultimately fruitless, particularly when dealing with dynamic

content (Chakrabarti et al., 1999). Focused crawls are far more efficient, but require an underlying understanding of the target networks' topology, including the topics of interest themselves.

Confusingly, terms such as 'dark web', 'deep web', 'invisible web' and 'hidden web' are often used interchangeably to denote a broad spectrum of somewhat exclusive concepts, with changing definitions including:

- Dynamically generated materials inaccessible via search engines due to the need for user input, rather than any desire for covertness or privacy (Florescu et al., 1998; Schadd et al., 2012);

* Corresponding author. Faculty of Information Technology, Monash University, 900 Princes Highway, Caulfield East, Victoria, Australia.

E-mail addresses: janis.dalins@monash.edu, janis.dalins@afp.gov.au (J. Dalins), campbell.wilson@monash.edu (C. Wilson), mark.carman@monash.edu (M. Carman).

- Unseemly or nefarious content such as that created by extremist/hate groups (Abbasi and Chen, 2007; Yang et al., 2009; Li et al., 2013); and
- Networks providing anonymity for content users and content providers (Iliou et al., 2016).

This paper exclusively refers to 'dark web' in a manner akin to (Iliou et al., 2016), though with further clarification. Following (Guitton, 2013), we regard anonymity to be the *non-coordinatability of traits* (Wallace, 1999). Non-coordinatability does not exist on the WWW, where any party involved in or capable of observing a particular communication can immediately glean each party's IP address, or at least that of an upstream provider¹. As alluded to by Guitton (2013), any sufficiently legally or technically empowered party could exploit such information in further efforts to identify the actual user(s). We regard dark webs to be networks which *at a technical level* do not rely upon elements capable of supporting coordinatability of traits - users can *choose* to make themselves identifiable, even inadvertently, but the network itself does not require them to do so as part of normal operation.

The scope of research into the 'dark web' within this work is restricted to Tor², a network built upon an openly published (and frequently revised) protocol employing temporary circuits of relays for the purposes of ensuring communicating parties' anonymity. Numerous software packages utilising Tor are freely available online, including a browser³. The browser provides a user experience entirely consistent with Mozilla Firefox⁴, which, if default settings are maintained⁵, provides even novice users with high level security and anonymity.

Akin to WWW domains and websites, Tor supports 'hidden sites', using an addressing system built upon randomly generated keys and denoted by an address ending in *.onion*, for example: <http://dxukgxbans5g5jn.onion>. Furthermore, Tor can be used to obtain anonymous access to WWW sites, meaning links to such sites are regularly observed on hidden sites.

The software's simplicity has arguably contributed to the popularity of Tor, with well-publicised services such as the *Silk Road* and *Agora* marketplaces effectively providing demand for users and types of commerce not traditionally associated with online retail trade.

The aim of this research is to investigate the structure of the Tor network and to develop a mechanism to efficiently identify locations likely to contain data of interest to law enforcement. In order for law enforcement to effectively search for content of interest, crawlers need to "understand" the content they are encountering, in so far as they should be able to efficiently classify the content according to its nature and likely value.

Related work

Identification of updated and changing domain names on the WWW is simple, with DNS providing a convenient, decentralised means for advertising resource addresses across the entirety of the publicly accessible internet. In contrast, hidden services issue *hidden service descriptors*. Unlike DNS, whereby a hierarchy of name servers takes responsibility for processing queries, hidden service descriptors are only published to a limited number of relay nodes within Tor. A requesting user queries the Tor network with a hidden

service's URL, receiving addressing information if the hidden service exists. Being decentralised and intentionally obfuscated, seeking and accessing hidden services is by its very nature slower than DNS. Ignoring the wider performance implications of a 'brute force' approach to identifying active services, such an approach would also be incredibly slow - there are 32^{16} possible top level domains⁶ for Tor hidden services, making a single threaded crawler experiencing extremely optimistic query times of 6 s per lookup take 2.3×10^{17} years to traverse each address.

Research into the topology of Tor using non-bootstrapped approaches is only possible through the use of technical exploits capable of bypassing the absence of obvious entry points. Biryukov et al. (2014), utilised an exploit available in February 2013 (since patched) to detect 39,824 hidden service addresses, their technical services offered, and their relative popularity. The sites' topics were classified using openly available software, though whilst the authors detail the software utilised, no details are given as to performance or accuracy. Reliance upon exploitation of bugs in Tor is unreliable at best, with the longevity of any associated crawls challenged by the fast pace of development and repair of existing Tor software - by way of example, the project's client software release notes⁷ - lists 150 versions, with over 600 minor and major bugfix summaries from underlying projects.

Guitton (2013) conducted a crawl of 1171 Tor hidden services, building a 23 category schema split broadly into *ethical* and *unethical* services - finding unethical services to be so pervasive that "... further development of Tor hidden services should hence stop".

Moore and Rid (2016) also conducted a crawl of Tor, creating a twelve category taxonomy. As with the other crawlers, the authors' aim was to provide a 'snapshot' of Tor and the services offered therein, without any specific interests or use cases. Of note, the authors specifically restricted their crawler to textual materials, due to the high risk of inadvertently accessing illegal materials such as child pornography and terrorist publications - a common issue within this field.

With respect to law enforcement applications, research into the dark web tends to focus upon specific content types. Chen (2012) establish a schema for categories of use for the web (also applied to 'dark webs'), focused upon terrorist organisations. Westlake et al. (2017), created a web crawler specifically focused upon locating child exploitation materials (CEM). Guided by seed websites plus a combination of keywords for identifying CEM plus 'safe sites' for domains previously identified as not being of interest, the crawler was successful in following the topic, though limited to three 'categories' based upon Canadian legal definitions of CEM. Whereas both studies ultimately relied upon manual labelling (with some automated features for steering and/or identifying content of interest), their focus upon particular content limits their use in identifying wider taxonomies.

Table 1 summarises and contrasts the various taxonomies developed in the aforementioned papers, ranging from particularly holistic (Biryukov et al., 2014) to targeted (Chen, 2012).

Automated classification of materials of interest is obviously of value in identifying illicit materials online, though as with previously detailed work, such work tends to be focused on narrow topics. For example, Fu et al. (2010), developed a crawler focused upon extremist discussion fora on the 'dark web'⁸, with a need to access and analyse multimedia driving in a mixed approach, with

¹ For example, a Virtual Private Network (VPN) host.

² The Onion Router.

³ 'Tor Browser' - refer <https://www.torproject.org/projects/torbrowser.html.en>.

⁴ Refer <https://www.mozilla.org/en-US/firefox/products/>.

⁵ Most notably, the disabling of Javascript.

⁶ Tor hidden service addresses consist of a sixteen character, base32 (a-z-2-7) string, appended with *.onion* - refer <https://www.torproject.org/docs/hidden-services.html.en>.

⁷ <https://git.torproject.org/tor.git>, as at 02 December 2017.

⁸ In this case, content requiring input in order to be accessible.

Download English Version:

<https://daneshyari.com/en/article/6884437>

Download Persian Version:

<https://daneshyari.com/article/6884437>

[Daneshyari.com](https://daneshyari.com)