

Understanding hidden websites deployed on Tor

Juha Nurmi

Tampere University of Technology

Email: juha.nurmi@ahmia.fi

Abstract—Tor is a software for anonymous TCP connections. This means that Tor enables anonymity to various Internet software. For instance, web servers can hide their location and web browsers can connect to these authenticated hidden services while the publisher and the viewer both stay anonymous. The publisher cannot be tracked down and the content cannot be censored. However, finding web content is laborious without an efficient search engine and therefore a search engine is needed for the Tor network. The aim of this paper is to introduce how to use our search engine implementation to understand hidden website

I. INTRODUCTION

Anonymity is an important right in order to support freedom of speech and defend human rights. An Internet user can use range of tools to hide ones identity[1]. Among these, the most popular tool is Tor. It has a large number various users, including ordinary citizens concerned about their privacy, corporations who do not want to reveal information to their competitors, and law enforcement and government intelligence agencies who need perform operations on the Internet without being noticed[2]. Further, human rights activist and journalist are communicating anonymously using Tor to protect their lives[3].

In addition, using the Tor network, it is possible to run web servers anonymously and without fear of censorship[4]. Servers configured to receive inbound connections through Tor are called hidden services (HSs); rather than revealing the real IP address of the server, a hidden service (HS) is accessed through the Tor network by mean of a virtual top level domain .onion[4].

In particular, we are interested in websites that operate as a hidden service. In this paper we call them hidden websites.

The published content is diverse[5]. Undoubtedly, some hidden websites are sharing pictures of child abuse, or operate as marketplaces for illegal drugs, including the widely known black market Silk Road. These few services are obviously controversial and often pointed out by critics of Tor and anonymity. On the other hand, vast number of hidden websites are devoted to human rights, freedom of speech, and information prohibited by oppressive governments.

Web search engines support finding web content. Because there were no search engines to search web content published using the Tor network, we built a working search engine for indexing, searching and cataloging content published inside the Tor network. Furthermore, we created an environment to share meaningful statistics, insights and news about the Tor network itself.

Ahmia provides the search and the access to hidden websites and believes that this is very important to the entire Tor network because we are efficiently enabling the diffusion and use of anonymous resources.

Whole search engine, Ahmia, is a free software and the source code is available online. This makes the research and our methods very transparent: Everyone is welcome to study our implementation.

In this paper we demonstrate how we can understand hidden service usage and how this reflects to our search engine design.

mds

December 04, 2014

II. RELATED WORK

The hidden wiki(s)

We are making HSs websites accessible in many ways. Because Ahmia is a public search engine for hidden services we would like them to be found and provide an easy gateway to visit hidden websites even without Tor software. To do this, we are supporting Tor2web Tor gateway project.

Tor2web is an HTTP proxy for Tor HSs designed initially by Aaron Swartz[6]. It aims at creating a network of proxies able to allow access to Tor HSs from the public internet. The software allows Tor hidden services to be reachable by means of a common browser and without the use of Tor client. Basically, it acts like a transparent proxy, translating the onion address into an HTTPS web URL[7].

In order to support tor2web.org, we maintain the Tor2web.fi proxy that enables people to connect to the .onion TLD with a regular web browser by replacing .onion part with .tor2web.fi. For example, <http://msydstlz2kzrdg.onion> can be accessed using <http://msydstlz2kzrdg.tor2web.fi>.

With the Tor2web software developers we introduced a HS discovery function to Tor2web software. This means that Tor2web is gathering a list of the visited .onion websites and the visit counts and search engines can download this list. In this way we can find new .onion domains and use their popularity information.

Ralf's paper (your [6]) would fit in there.

Maybe even say that Metrics, while having lots of statistics on the Tor network, doesn't have anything good about hidden services and their content.

III. FINDING, RANKING AND UNDERSTANDING CONTENT

First, the start point is to crawl the web content form the hidden services. Before that can be performed a seed list of .onion domains is used. However, Tor technology does

Unfortunately, this method finds only those new .onion sites which are linked to those .onion pages which are already indexed. Moreover, only few hidden websites links to other hidden websites and there is no linking to every .onion. As a result, we cannot find all hidden sites. We visualized this problem by generating a SVG image of the crawling paths (figure 1). More visualizations material is available on <https://ahmia.fi/static/visuals/>.

We would like to show a glance of the hidden website content in general. Ahmia produced a word cloud visualization of the front pages of hidden websites (see figure 2).

Similarly, using the search index of hidden websites, we are locating malicious software sites to inform security firms, child pornography sites to filter them out, and immediately after the international law enforcement operation, Operation Onymous, the list of sites seized by them.

- [1] Goldschlag, David, Michael Reed, and Paul Syverson. "Onion routing." *Communications of the ACM* 42.2 (1999): 39-41.
- [2] Dingledine, Roger, Nick Mathewson, and Paul Syverson. "Deploying low-latency anonymity: Design challenges and social factors." *Security & Privacy, IEEE* 5.5 (2007): 83-87.
- [3] Tor Project homepage. The Tor Project, Inc. <https://www.torproject.org/>
- [4] Dingledine, Roger, Nick Mathewson, and Paul Syverson. *Tor: The second-generation onion router*. Naval Research Lab Washington DC, 2004.
- [5] Biryukov, Alex, Ivan Pustogarov, and Ralf-Philipp Weinmann. "Content and popularity analysis of Tor hidden services." *arXiv preprint arXiv:1308.6768* (2013).
- [6] <http://www.aaronsw.com/weblog/tor2web>
- [7] <http://logjoshermes.org/home/projects-technologies/tor2web/>

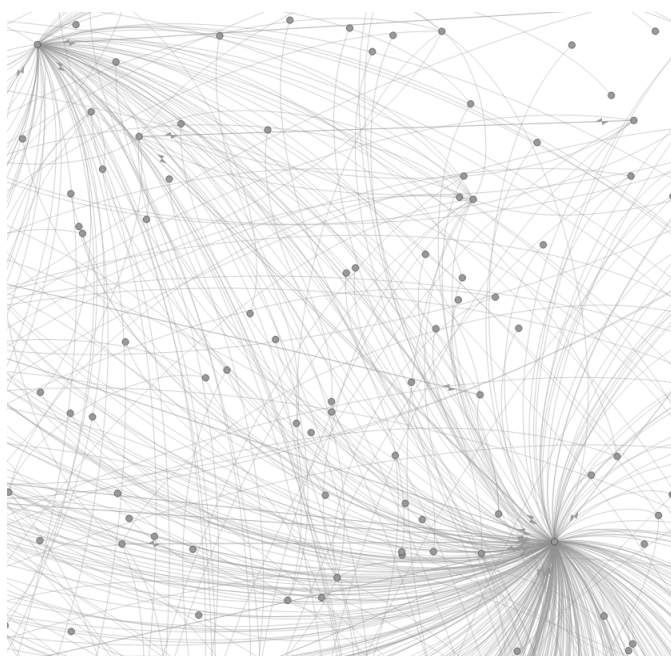


Fig. 1. A part of the visualization of the linking structure of hidden websites. Few sites gather lists of other .onion domains while the most of the sites have no linking to other .onion domains.



Fig. 2. A word cloud that represents the most popular text content. These are the most common words on the front pages of hidden websites.

By Tor2web average visits	By public WWW backlinks
216003 pinkmethuylnenlz.onion	24538 strngbxhwyuu37a3.onion
22873 t54cjs4qc2r4bn63.onion	3252 kpzv7ki2v5agwt35.onion
13223 3qwajq5p5pfjs3sw.onion	2852 silkroad6ownowfk.onion
13064 h3vf5lellsvjlqlx.onion	1830 silkroad5v7dywlc.onion
12799 torbookd7whjnj4u.onion	1520 3g2upl4pq6kufc4m.onion
12239 svcz25e3m4mwlauz.onion	1510 am4wuhz3zifexz5u.onion
11414 npdaaf3s3f2xrmlo.onion	1410 grams7enufi7jmdl.onion
7756 64ansq6xm5mmsb3a.onion	1350 xmh57jrznw6lnsl.onion
7266 juvatztgkarpzp2o.onion	1034 silkroadvb5piz3r.onion
6530 girlshitlrelazwm.onion	1020 agorahooawayvfoe.onion

Fig. 3. The most popular websites according to average Tor2web proxy visits per day and the most popular websites according to number of backlinks from the public WWW.