

Climbing China's Great Firewall

Adam L. Casey
Division of Science and Mathematics
University of Minnesota, Morris
Morris, Minnesota, USA 56267
casey369@morris.umn.edu

ABSTRACT

Many different countries censor the internet within their state. Citizens frequently wish to avoid the state censorship. There are many different methods that have been developed to achieve this. Governments and citizens are in a constant arms race, with both developing opposing technologies. China in particular has the largest population of people on the planet, and the Chinese government attempts to censor the internet. This paper will investigate three methods of navigating around state censorship: Cachebrowser, INTANG and Tor. Cachebrowser and INTANG were developed specifically to navigate around state censorship while Tor was originally developed for anonymous browsing. This paper will analyze their effectiveness and viability to avoid censorship.

Keywords

China, Great Firewall, Tor, CDNs, INTANG

1. INTRODUCTION

The largest group of people in the world with censored access to the internet is the population of China, with upwards of 1.3 billion people. The Chinese government still actively censors the internet through a system of network monitoring and network manipulation, referred to broadly as the Great Firewall of China (GFW) [11].

This paper will describe the methods and evaluate the success of three separate ways of evading the censorship of the GFW. These methods are Cachebrowser, INTANG and Tor. These methods each work in very different ways. Cachebrowser works by allowing the user to easily access uncensored versions of websites that are usually blocked by accessing the cached versions of these sites that the Chinese government cannot feasibly block for reasons that will be discussed in section 3.1. INTANG works by manipulating the internet traffic being sent by the user's machine directly to avoid censorship by manipulating packets which will be discussed in depth in section 3.2. Tor works by bouncing the user's connection through multiple different nodes. This makes it difficult to track. How the Chinese government probes for Tor servers will be discussed in section 3.3

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/4.0/>.
UMM CSci Senior Seminar Conference, April 2018 Morris, MN.

2. BACKGROUND

In order to understand how circumvention techniques work a basic understanding of some internet protocols and concepts is required. In this section, background will first be given on the basic frameworks of the internet. Also, background is given on the Transport Control Protocol (TCP). This background information is necessary to understand how INTANG circumvents the GFW. Information is also given on Content Delivery Networks (CDNs). A foundation in this subject is necessary to understand how Cachebrowser circumvents the GFW. An overview of Tor will also be given in this section. Background on Tor is necessary to understand one of the most common ways of circumventing internet censorship and to understand how the Chinese government attempts to stop Tor traffic.

2.1 Internet Basics

In order to understand how the Chinese government censors the internet and to understand the TCP protocol, first some internet frameworks must be explained. To begin with client and server roles need to be explained. When you make a request to visit a website such as www.facebook.com your computer talks to the server for facebook. This is referred to as the server because the user makes a request and the server serves the content to the user. The user in this case is referred to as the client.

Before the client actually makes a request to facebook.com it needs to know how to get to it. Each server on the internet has a unique address associated with it. This is called an IP address [10]. In order for the client to find the IP address for a server they are trying to access they must make a Domain Name System (DNS) request. This is done by making a request to a DNS server which has a list of IP addresses for different websites in a cache. The client tells the DNS server the website they are trying to access. The server then sends the client the IP address of the website they are trying to access. Finally, the client can make a request to the server they are trying to access and send data to it and the server can send data back to the client.

One way the Chinese government attempts to censor the internet is by manipulating the records on DNS servers to return incorrect IPs for websites or to drop packets for requests to websites the government does not want clients to be able to access. This practice is referred to broadly as DNS poisoning [3].

Another way the Chinese government attempts to censor the internet is through a technique called IP address filtering [5]. This is where the IP of the website you are trying to

TCP Header																																	
Offsets	Octet	0								1								2								3							
Octet	Bit	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
0	0	Source port																Destination port															
4	32	Sequence number																															
8	64	Acknowledgment number (if ACK set)																															
12	96	Data offset	Reserved 0 0 0			N S	C W R	E C E	U R G	A C K	P R H	R S T	S Y N	F I N	Window Size																		
16	128	Checksum																Urgent pointer (if URG set)															
20	160	Options (if data offset > 5. Padded at the end with "0" bytes if necessary.)																															
...																															

Figure 1: Diagram of a TCP packet taken from [9].

access is blocked at the hardware level, usually on your home router. This is different from DNS poisoning because even if we know the IP address of the website we are attempting to access we cannot get to it.

The final censorship technique that will be mentioned in this paper is keyword censorship [11]. This is where the Chinese government monitors the client's internet traffic and will terminate connections with keywords the government has deemed sensitive. This is a way to block connections if the IP of the website is not blocked and its DNS records have not been poisoned.

Another concept that is necessary to understand circumvention techniques is the idea of a proxy, sometimes referred to as a proxy server [7]. A proxy is essentially a server that makes web requests for you. Instead of talking directly to the server the client wants to access, it talks to a different server, which talks to the server the client wants to talk to. Then the proxy server then takes the information the intended server sent it and sends it back to the client.

2.2 The TCP Protocol

Transport Control Protocol (TCP) [10] is one of the main internet standards. The main data structure of TCP is the packet. Data is sent in discrete pieces in order. One of these pieces is referred to as packet. This order is maintained by incrementing the sequence number with each new packet that is sent. Packets are also sent to establish and end a connection. Figure 1 shows a diagram of a TCP packet. The first 16 bits are reserved for the source port. The next 16 are reserved for the destination port. Bits 32 to 63 are reserved for the sequence number. Bits 64 to 95 are reserved for the acknowledgement number if the packet is flagged as an ACK (acknowledgement) packet. The data offset bits serve a dual purpose. They indicate the size of the TCP header and the offset from the beginning of the header to the beginning of the actual data. The next three bits are reserved. Next, the flag bits begin. These indicate the type of packet. Bits 112 to 128 indicate the windows size the sender of the packet is willing to receive. The next section of bits is a checksum which is used for error checking. The next section is the urgent data if this is a flagged URG (urgent) packet. There are multiple flags that can be set. Many of these are not relevant

to this paper the relevant ones are SYN/ACK/FIN/RST. The ACK flag indicates that the packet is an acknowledgement packet. The RST flag indicates that the connection should be reset. The SYN flag indicates that sequence numbers should be synced. In practice this should only be used at the beginning of the connection process. The FIN (finish) flag indicates that this is the last packet from the sender.

Another important aspect of the TCP protocol to know for this paper is Time to live (TTL). Each packet is given a TTL. The TTL indicates how long the packet will stay in the network before it destroys itself. This is useful so that packets do not clog up network infrastructure forever if they do not reach their intended destination.

There are three main steps to transmitting data using the TCP protocol. This first step is to establish the connection using a handshake process. A connection is a link between client and host where data can be sent freely between client and server. Next the actual data is sent. Finally the connection is terminated.

Connection establishment is a multistage process. The first stage is when the user establishes a connection to the server. Next the user sends a SYN (sync) packet to the server. Once the server receives this SYN packet it sends an ACK (Acknowledgement) packet back to the user along with a SYN packet. The next step to establish data exchange is for the user to send an ACK packet back to the server. After this process has been completed regular data transfer can occur.

Closing a connection is also a multistage process. Closing a connection is different from establishing a connection in that a client or a server can close the connection. Only a client can establish a connection. The process for closing a connection on the server end or the client end is the same but the roles are reversed. As an example let's consider a server closing a connection with a client. The server will send a FIN packet to the client with the sequence number of the next data packet the client is expecting to receive. Once the client receives the FIN packet from the sever it sends an ACK packet with the sequence number increased by one. The client then sends the server its own FIN packet with sequence number relative to the amount of data it has sent to server thus far. The server will then send a final ACK

packet to the client and the connection will be terminated on both ends.

Another part of the TCP protocol that is manipulated to get around censorship is the TCP Control Block (TCB). The TCB is a data structure that is created by the TCP protocol when a connection is established. This TCB is normally created on the client and the server, but in the case of the GFW one is also created by the GFW to monitor the connection. The purpose of the TCB is to keep track of all connections incoming and outgoing on the machine it is created on. The GFW uses the TCB it creates in combination with packet inspection to terminate connections with sensitive keywords.

2.3 CDNs and cached content

CDNs [6] also known as content delivery networks are a distributed set of servers hosting web content. The goal of this system is to decrease latency by redirecting users to servers hosting the content near them instead of one central one that could possibly be across the globe. CDNs are not run by the companies that have the actual content. They are run by a separate company and pay the company that runs the CDN to host their content. There are many reasons to use a CDN. One reason is less stress on a single server. Since the network load is spread out between multiple servers hosting the cached content, no one server takes the brunt of the load. The servers that host this cached content that users access are referred to as edge servers.

A common way to implement a CDN for an already existing website is through DNS modifications. When navigating to the web address for a site, the client will be redirected for the CDN server for all content that does not frequently change on the site. This would be things like logos, headers that are always the same, etc. Content that changes frequently will be served by the host server. In some cases for very popular websites dynamic content is still hosted on CDN servers. This is accomplished by having a high speed pipeline from the host server to the CDN server which keeps the CDN server up to date. When a CDN server has content change on it, it relays this information to all other CDN servers hosting the site so that things are consistent. As a result of the fact that CDNs are operated by companies separate from the ones wanting their web content hosted, multiple different websites content is stored on the same CDN server. The client just requests the portion of the content that they actually need.

2.4 Tor

Tor [8] gets its name from the project's previous name "The Onion Router". Tor is a free piece of software that is used for anonymous internet browsing. It achieves this by redirecting the user's connection through multiple different nodes. Tor is referred to as "The Onion Router" because each data packet is wrapped in a layer of encryption for each node that it passes through. Each node only decrypts one layer of information to know where to send the packet next and does not have access to the actual data you are sending. This is because the actual data is only located in the last layer. Figure 2 illustrates this process of traffic traveling through multiple nodes. This makes it difficult to track. Another important feature of Tor to understand is the Tor bridge. A Tor bridge is essentially the same as a Tor node. The difference is that the list of Tor bridges is

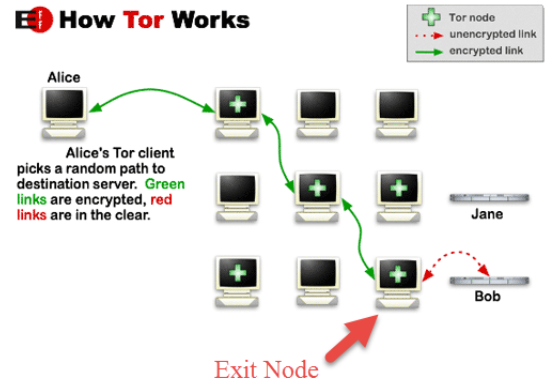


Figure 2: Diagram of Tor traffic through nodes. Taken from [1].

not publicly available, like it is for Tor nodes. One of the ways the Chinese government blocks Tor traffic is to block all IPs of the publicly listed Tor nodes.

Tor also has the ability to use multiple different pluggable transports. A pluggable transport is a type of cipher suite, a set of encryption algorithms which the data is encrypted with before being sent out. This allows users to access parts of the Tor network that are usually blocked by the Chinese government using IP address filtering. Pluggable transports also work to disguise traffic from being identifiable as Tor traffic. This is because even if censors do not know what website you are trying to access they will terminate your connection if they recognize it as Tor. Not all pluggable transports work well for avoiding the censorship of the Great Firewall however. These pluggable transports often have only one layer of encryption and encrypt packets in an easily recognizable way, allowing the government to recognize Tor traffic and block it accordingly. The currently recommended pluggable transport to use is meek-amazon. Meek-amazon works by using a technique called domain fronting. While too complicated to explain in-depth in this paper, domain fronting makes it look like the client is accessing a different website than they actually are. The meek-amazon pluggable transport accomplishes this by routing traffic through Amazon cloud servers, which then access the Tor network as a proxy. Figure 3 illustrates how the meek-amazon pluggable transport works.

3. METHODS OF CIRCUMVENTION

3.1 Cachebrowser

Cachebrowser [11] is a tool developed by John Holowczak and Amir Houmansadr to bypass the censorship of the Great Firewall of China. It uses CDNs and cached content, as mentioned in Section 2.3, to access web pages that would be inaccessible using the internet without a circumvention tool. Using cached content is way of circumventing common censorship techniques such as IP address filtering. IP address filtering is the process of blocking access to certain IP addresses. IP address filtering is ineffective at blocking CDNs because one website is spread across multiple different IPs, so the censors would have to blacklist all of them to block the site. Also, because of the way edge servers work, multi-

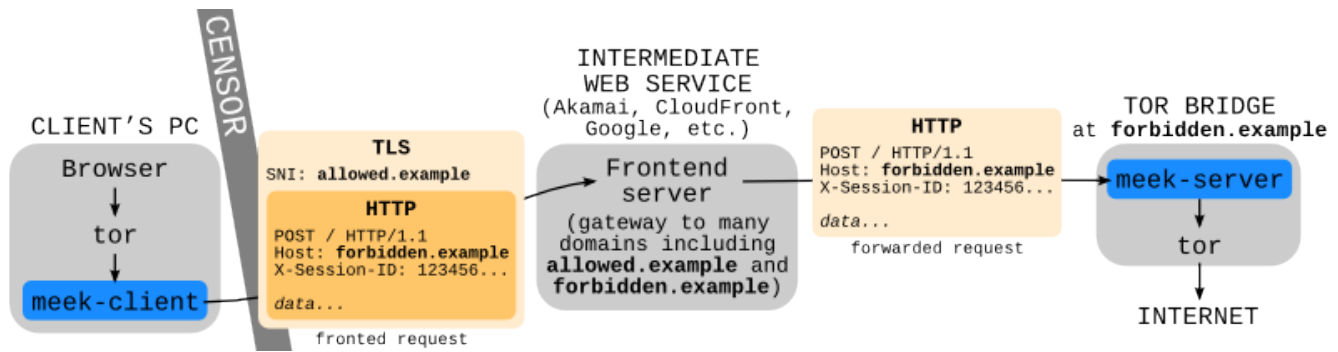


Figure 3: Diagram of meek-amazon domain fronting taken from [4].

ple sites are hosted on one server, so blocking access to one server will inadvertently block access to all websites on that shared IP, not all of which the censor necessarily desires to block.

DNS interference is another widely used censorship technique and is most effective at blocking cached content. DNS interference works by interfering with the name resolution process when trying to access a blocked website. This is effective at blocking cached content because it doesn't matter how many IPs the content is spread across or if more than one site is hosted at that one IP. This is because the DNS interference prevents the end user from knowing the IP of the content in the first place.

Cachebrowser works by keeping its own database of IP addresses to CDN hosted alternatives to regular content. When Cachebrowser encounters a CDN domain it internally enumerates and saves all other IP addresses for that CDN in case that one is censored. If Cachebrowser encounters a customer domain it will return the addresses of CDNs which host that content. This is done by using the free DNS resolver www.digwebinterface.com. This site is currently not blocked in China.

As an alternative to this method, Cachebrowser also implements a remote bootstrapper using SWEET [2]. SWEET is a communications tool that encapsulates messages through emails and sends them through standard email protocols. In the case of Cachebrowser, a web server is set up in the United States watching for emails. When a Cachebrowser user makes a request for a site it does not have a CDN-hosted IP for a message is sent out through SWEET using the bootstrapper. This message is sent to a server in the United States. The server then makes the DNS lookup and sends the information back to the client to be added to its database.

3.2 INTANG

INTANG [5] is a tool developed by Wang et al. to avoid the censorship of the GFW. It combines many different strategies. A large portion of the paper [11] is dedicated to figuring out how the GFW works using trial and error to gain knowledge of the system in order to help develop a tool to avoid it. INTANG works by implementing all of strategies evaluated in the analysis section of the paper. There are too many individual strategies to go over all of them in this paper, but they fall into 3 main categories: TCB creation, TCB teardown, and data reassembly.

TCB creation works by sending a SYN insertion packet

with the incorrect sequence number to create a false TCB on the GFW and then initiating the real connection with the server, which the GFW will ignore because of the sequence number discrepancies.

TCB teardown works by by crafting RST,RST/ACK and FIN packets with a TTL constructed in a such that the packet reaches the GFW and terminates the TCB but does not reach the server, thus keeping the server alive.

Data reassembly has two separate forms: Out-of-order data overlapping and In-order data overlapping. Out-of-order data overlapping works by sending garbage data fragments with the same offset and length as the real data. When the GFW encounters two packets with the same offset and length it records the first one and ignores the second. In-order data overlapping works by filling up the GFWs input buffer until it is overloaded and can no longer read new data. This done by crafting insertion packets with either a wrong checksum or a very short Time to Live (TTL) so they fill up the GFW buffer while still keeping the connection to the server alive.

INTANG uses a combination of all these strategies together with new strategies developed from analyzing the performance of the old strategies to approach circumvention from multiple angles. INTANG is a measurement tool, which means that it keeps track of which strategies work with regards to different IPs and adjusts strategies in use based on that.

3.3 Tor

Unlike [5] and [11] which both developed tools to circumvent the GFW and tested them within their paper, [4] instead examines the GFW behavior and hypothesizes on how this information can be used create Tor servers that can more easily avoid blacklisting. Probing for Tor servers is triggered the the GFW sees traffic that carries the signature of a cipher suite that Tor uses. Tor in its default configuration is essentially completely useless already since it has publicly available list of IPs the government can blacklist outright. Tor developed a pluggable transport layer in response. The suggested pluggable transport to use currently is meek-amazon or meek-azure. Neither of these transports were tested in the paper as they were new and not yet popular at the time.

For the study, two different infrastructures were created and long running Tor servers' logs were analyzed. One infrastructure that was created was the shadow infrastructure. This infrastructure was built up of private Tor bridges that

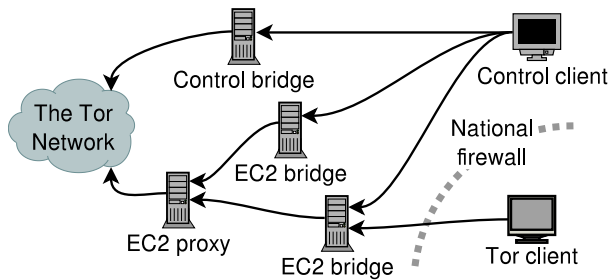


Figure 4: Diagram of the Shadow dataset taken from [4].

DNS resolver	IP	except Tianjin	All
DYN 1	216.146.35.35	98.6 %	92.7 %
DYN 2	216.146.36.36	99.6 %	93.1 %

Figure 5: DNS server access success rates. Taken from [11].

only the researchers could access. Figure 4 is diagram of how the Tor clients and servers were set up for the shadow data set.

The second infrastructure was the Sybil infrastructure. The Sybil infrastructure set up a Tor bridge in France that redirected 600 ports to it using a firewall redirection, then connected to each of the ports in ascending order. The Log dataset came from analyzing the logs a one of the researcher’s Tor servers that has been running since January 2011. By looking at logs they found the server had been subject to active probing from China for 2.5 years, first showing up in 2013. This was not the result of an attempt to induce probing but seemingly the regular amount of probing by the government. The logs were then used to evaluate how effective probing was at disrupting Tor by looking at whether the TCP handshake was completed. They found that obfs2 and obfs3, which are different types of pluggable transports, were very rarely disrupted and had high success rates while Tor without a pluggable transport was essentially unusable.

4. RESULTS

4.1 Cachebrowser Results

Cachebrowser [5] was able to successfully load facebook.com on a client’s machine. This is a complex website with content hosted on multiple CDNs that is completely inaccessible to users not using some kind of circumvention technique. Research found that of the top 1000 Alexa websites, 82% were hosted on some CDN provider while 85% of news sites like wsj.com were hosted on CDNs.

The main issue with Cachebrowser is latency. Figure 6 compares the latency times of Cachebrowser and other alternative methods of circumvention within and outside of China. As you may notice facebook.com does not have a latency value for inside of China. This is because facebook.com is completely censored within China. Also, there is no latency measurement for Tor within China. The researchers indicate that they were not allowed to run Tor on their client in China and that is why it is not included here. Cachebrowser does have higher latency than non-

censored browsing in every case where non-censored browsing is available. The maximum difference between page download times between Cachebrowser and a non-censored version inside of China is .725 seconds. Moving on to the latency sample data from page access in the United States, Tor has higher latency than Cachebrowser in every case. This suggests that if Tor were to be tested within China it would most likely also be slower than Cachebrowser there as well. The privacy of Cachebrowser is also robust. Assuming the CDN is using an encrypted pathway, which almost all do, the state cannot see what content you are viewing. The state will also not know what website in particular you are viewing because multiple sites are hosted at the IP since the CDN serves multiple different sites off of the same server. The only possible leak in this situation is that the CDN itself has information about its visitors. However, sharing this information with anyone is a violation of the CDNs user agreement therefore we can assume that they do not share this information.

4.2 INTANG Results

INTANG [11] in general is very successful at avoiding the censorship of the GFW. Research indicates that success rates for different strategies vary wildly. INTANG takes advantage of this by implementing multiple different packet manipulation strategies and choosing the best one systematically. Figure 7 shows the success rates of the multiple individual strategies tested and the success rate of INTANG itself.

One side effect of INTANG is that because the government poisons DNS requests in the same way it blocks TCP traffic, INTANG is also effective at evading DNS poisoning. Figure 5 shows the success rate of accessing DNS servers using INTANG. The data was collected querying a DNS poisoned domain, in this case www.dropbox.com, 100 times. The discrepancy in the DNS resolution table is because the area of Tianjin is an outlier. Success rates are much lower in Tianjin. The exact reason for this is unknown, but it is hypothesized that the GFW infrastructure is more developed in this area. Success rates in Tianjin were 38% for DYN 1 and 24% for DYN 2. DYN 1 and DYN 2 are two different DNS servers. In particular DYN 1 is Google’s DNS server 8.8.8.8 and DYN 2 is Google’s DNS server 8.8.4.4. Normally connections to these DNS servers are terminated using a TCP connection termination. Research also indicates that two OpenDNS resolvers 208.67.222.222 and 208.67.220.220 are uncensored even without the use of INTANG. This is similar to the method that the bootstrapper uses of simply request the DNS lookup from an uncensored address in the first place.

4.3 Tor Results

Research [4] found that while the GFW may not be able to detect Tor traffic from clients using sufficiently advanced pluggable transports such as obfs2 and obfs3, it can effectively probe for and shut down Tor servers regardless of the cipher suite used. Research also indicated that the Chinese government’s Tor probes are active in real-time, only stopping for short periods.

5. CONCLUSIONS

After analyzing three circumvention techniques it is clear that there are advantages and downsides to each. Tor is

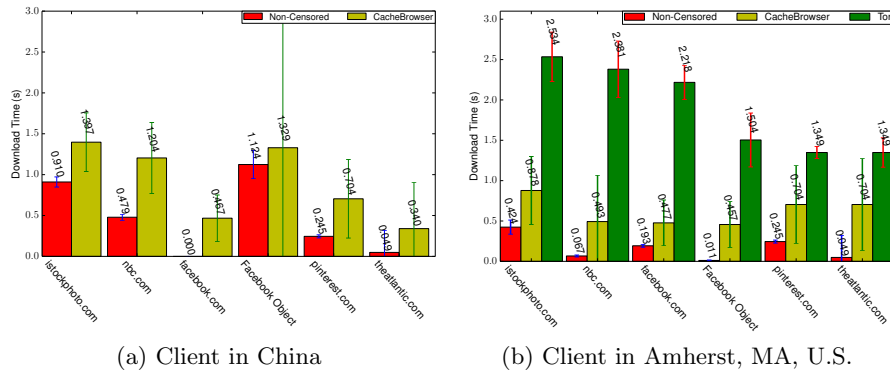


Figure 6: Cachebrowser latency compared to regular access latency and Tor latency inside China and outside of China. Taken from [5].

Vantage Points	Strategy	Success		
		Min	Max	Avg.
Inside China	Improved TCB Teardown	89.2%	98.2%	95.8%
	Improved In-order Data Overlapping	86.7%	97.1%	94.5%
	TCB Creation + Resync/Desync	88.5%	98.1%	95.6%
	TCB Teardown + TCB Reversal	90.2%	98.2%	96.2%
	INTANG Performance	93.7%	100.0%	98.3%

Figure 7: INTANG success rates taken from [11].

proven to be effective and has an active development community. It is, however, slow and prone to frequent shut-downs. This is due to the government identifying Tor servers through probing. Cachebrowser is an effective way of viewing websites that would normally be censored by the Chinese government, but it has downsides as well. If a website is not hosted on a CDN it cannot be accessed using Cachebrowser at all. INTANG is useful for navigating around keyword censorship, however it is prone to updates to the GFW. If the government implements changes that negate the packet manipulations done by INTANG the entire tool is rendered useless. I would suggest using a combination of INTANG and Cachebrowser to avoid censorship first, as they have lower latency than Tor. If the content you are trying to view cannot be accessed using this strategy, Tor can be used as a backup.

Acknowledgments

I would like to thank Elena Machkasova and Nic McPhee for helpful feedback on drafts of this paper. I would also like to thank UMM alumni Jeff Lindblom for their useful feedback as well.

References

[1] Nandan Desai. *Unclosing the Dark Web - Post 2*. June 2016. URL: <https://knowledgetransmitter.quora.com/Unclosing-the-Dark-Web-Post-2>.

[2] G. Dlodla, M. Bembe, T. Olwal, and Jun Kyun Choi. “Enhanced SWEET protocol for energy efficient wireless sensor networks”. In: *2013 International Conference on ICT Convergence (ICTC)*. Oct. 2013, pp. 332–335.

[3] *DNS spoofing*. Apr. 2018. URL: https://en.wikipedia.org/wiki/DNS_spoofing.

[4] Roya Ensafi, David Fifield, Philipp Winter, Nick Fearnster, Nicholas Weaver, and Vern Paxson. “Examining How the Great Firewall Discovers Hidden Circumvention Servers”. In: *Proceedings of the 2015 Internet Measurement Conference*. IMC ’15. Tokyo, Japan: ACM, 2015, pp. 445–458. ISBN: 978-1-4503-3848-6.

[5] John Holowczak and Amir Houmansadr. “CacheBrowser: Bypassing Chinese Censorship Without Proxies Using Cached Content”. In: *Proceedings of the 22Nd ACM SIGSAC Conference on Computer and Communications Security*. CCS ’15. Denver, Colorado, USA: ACM, 2015, pp. 70–83. ISBN: 978-1-4503-3832-5.

[6] *How Content Delivery Networks Work*. URL: <https://www.cdnetworks.com/en/news/how-content-delivery-networks-work/4258>.

[7] *Proxy server*. Apr. 2018. URL: https://en.wikipedia.org/wiki/Proxy_server.

[8] *Tor*. URL: <https://www.torproject.org/about/overview.html.en>.

[9] *Transmission Control Protocol*. Mar. 2018. URL: https://en.wikipedia.org/wiki/Transmission_Control_Protocol.

[10] *Transmission Control Protocol*. URL: <https://tools.ietf.org/html/rfc793>.

[11] Zhongjie Wang, Yue Cao, Zhiyun Qian, Chengyu Song, and Srikanth V. Krishnamurthy. “Your State is Not Mine: A Closer Look at Evading Stateful Internet Censorship”. In: *Proceedings of the 2017 Internet Measurement Conference*. IMC ’17. London, United Kingdom: ACM, 2017, pp. 114–127. ISBN: 978-1-4503-5118-8.