

# Systemization of Pluggable Transports for Censorship Resistance

Sheharbano Khattak\*, Laurent Simon\*, Steven J. Murdoch†

\*Computer Laboratory, University of Cambridge, UK

forename.lastname@cl.cam.ac.uk

†Computer Science Department, University College London, UK

s.murdoch@ucl.ac.uk

**Abstract**—An increasing number of countries implement Internet censorship at different levels and for a variety of reasons. The link between the censored client and entry point to the uncensored communication system is a frequent target of censorship due to the ease with which a nation-state censor can control this. The diversity of a censor’s attack landscape has led to an arms race, leading to a dramatic speed of evolution of censorship resistance schemes (CRSs) (we note that at least six CRSs have been written in 2014 so far). Despite the inherent complexity of CRSs and the breadth of work in this area, there is no principled way to evaluate individual systems and compare them against each other.

In this paper, we (i) sketch an attack model to comprehensively explore a censor’s capabilities, (ii) present an abstract model of a *Pluggable Transport* (PT)—a system that helps a censored client communicate with a server over the Internet while resisting censorship, (iii) describe an evaluation stack that presents a layered approach to evaluate PT, and (iv) survey 34 existing PTs and present a detailed evaluation of 6 of these corresponding to our attack model and evaluation framework. We highlight the inflexibility of current PTs to lend themselves to feature sharability for broader defense coverage. To address this, we present *Tweakable Transports*—PTs built out of re-usable components following the evaluation stack architecture with a view to flexibly combine complementary PT features. We also list a set of challenges to guide future work on *Tweakable Transports*.

## I. INTRODUCTION

As the Internet becomes an increasingly important means to engage in civil society, those who wish to control the flow of information are turning to measures to suppress speech which is considered undesirable. While blocking can take place at any point(s) in the network, the link between the censored client and entry point to the uncensored communication system has been a frequent target<sup>1</sup>. This is so because the censor is typically a powerful nation-state adversary, and has control over network infrastructure within the censored region. Consequently there is a growing demand for Censorship Resistance Schemes (CRSs) which can bypass these blocks. A comprehensive CRS must defend against all blocking techniques available to censors, and so the schemes have become increasingly complex. As no one scheme has proved resistant to all potential adversaries, an arms race has developed resulting in the evolution of blocking resistance techniques to have dramatically sped up. This is well captured in Figure 1 on page 1 which shows work in this area over the last five years as per our survey.

The diversity of a censor’s attack landscape and the profusion of CRSs that defend against different attack paths makes

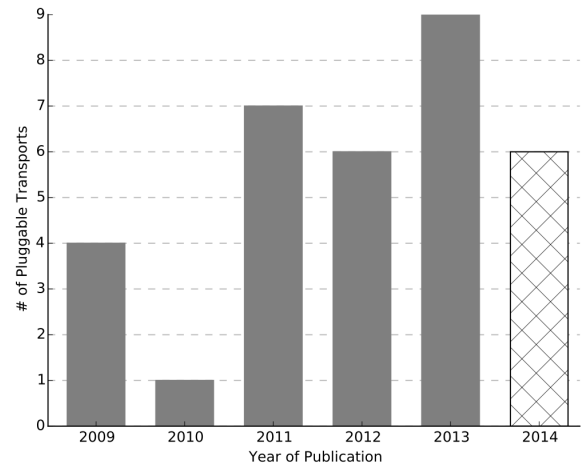


Fig. 1: Surveyed systems (literature and implementation) from the last five years that concern link obfuscation.

it hard to evaluate individual tools<sup>2</sup> and compare them against each other to identify gaps. In this paper, we sketch a comprehensive attack model to understand a censor’s capabilities and the circumvention scope of various CRSs. Next we propose an abstract model of a system—a *Pluggable Transport* (PT)—that enables a client application in censored region to communicate with a server application over the Internet, even though direct connections are blocked.<sup>3</sup> We then outline an evaluation stack that can be used to understand capabilities of various PTs in terms of the attack path(s) these protect. In this paper, we survey 34 papers and map them to three high level classes according to the path on attack diagram these seek to protect. For each class, we evaluate the significant work in that area in a systematic way according to our attack model and evaluation stack.

We note that combining complementary PT features provides broader defense on censor’s attack model, however this is not practical because most PTs have been designed as monolithic systems. There has been some effort to chain *Pluggable Transports* in a blackbox fashion [65] or by adaptation of source code. This approach however does not offer seamless integration and suffers from temporal overhead: fast

<sup>2</sup>Adversary Lab [5] has done some preliminary work on how to evaluate individual tools by running them in a standard environment.

<sup>3</sup>The concept of PTs is not entirely novel and has been the de facto API for anonymous communication systems to integrate with censorship resistance schemes [69], however, we purposely maintain a broad focus to accommodate CRSs that have not been written strictly as PT but can fit the broad model with minor adaptation.

<sup>1</sup>Through the rest of this paper, we limit our scope to this link.

development is particularly important for censorship resistance because there is no one approach which is optimally efficient and resistant to all attackers. Given these issues, we propose an extension of Pluggable Transports—Tweakable Transports. Tweakable Transports are Pluggable Transports built out of re-usable components following the evaluation stack architecture. Each component can be replaced with another which is compatible and components can be inserted or removed. This approach allows code-reuse because a component developed for one Tweakable Transport can be used for another. In so doing, more collaboration opportunities are allowed, better testing can be performed on frequently required components improving reliability and both spatial and temporal agility. As a result Tweakable Transports exponentially increase the number of possible CRSs.

The contributions of this paper are as follows:

- We present an attack model from a censor’s perspective that captures the diversity of censorship mechanisms and models the capabilities of a censor (Section II).
- We present an abstract model of a Pluggable Transport—a system that facilitates communication between a censored client and server in a censorship resistant fashion (Section III). Our model is derived from the existing protocol specification [69], but is generic enough that most CRSs not originally written as PTs can fit it with minor modifications.
- We outline an evaluation stack for representing PT capabilities (Section III). This stack is based on functional capabilities of a PT as per attack paths these seek to protect.
- We survey 34 PTs and identify broad classes with respect to the attack path(s) that these seek to protect (Sections IV, V and VI). Furthermore, we use our attack model and evaluation stack to comprehensively understand the assumed threat landscape and defenses offered by these PTs. We find that most PTs tend to cluster either around content filtering resistance or IP filtering resistance. This is an unrealistic assumption because most censors are capable of performing both kinds of censorship.
- We note that a comprehensive circumvention scheme can benefit from integration of complementary PT features. This is, however, not possible as most PTs have been designed in a monolithic fashion. Motivated by this observation, we propose *Tweakable Transports*—PTs built out of re-usable components following the evaluation stack architecture (Section VII).

## II. A CENSOR’S ATTACK MODEL

The circumvention technologies discussed in this paper will be evaluated against the censor’s attack model illustrated in Figure 2 on page 3. The overall goal of an attacker is to disrupt access to certain material, while minimising disruption to other material and to do so the targeting may be direct or indirect. With direct targeting, there is a static feature of traffic that efficiently and accurately identifies that the traffic should be disrupted (e.g. by matching the IP address against

an accurate blacklist). Indirect targeting is used when direct targeting is not possible, i.e. when no suitable static feature exists. In this case traffic is first fingerprinted (e.g. based on timing characteristics) so as to derive a set of features (e.g. IP addresses to which traffic matches the fingerprint) which can then be used for direct targeting.

Once a decision to disrupt traffic has been made, a censor can either corrupt it by inserting false information, by deleting and/or modifying existing information or by disabling its access and/or distribution<sup>4</sup>. In Figure 2 on page 3, this is represented using green and blue color codes. A censor can realize these goals on any (combination) of points along the information dissemination infrastructure: client side—the information consumer, server side—the information server, or the channel over which information travels.

### A. Blocking

#### On Client System (CEN.CLI).

A censor can directly enforce blocking on the client-side system, for example by installing surveillance software openly, or discretely by compromising the system (for example by means of malware or insider attacks). Once a censor has gained the desired control over client system, there is a range of things it can do, including but not limited to corruption of incoming network traffic before it reaches the application program, scanning data for blacklisted keywords and communicating results to an upstream server for subsequent blocking, and blocking outgoing connections to a blacklist of destination IP addresses. Effectively, a censor can disable access as well as corrupt information on the client system. China’s Green Dam, a filtering software product purported to prevent children from harmful Internet content, was mandated to be installed on all new Chinese computers in 2009 [29]. The software was found to be far more intrusive than officially portrayed, blocking access to a large blacklist of websites in diverse categories, and monitored and disrupted operation of various programs if found to be engaging in censored activity. TOM-Skype, a joint venture between a Chinese telephony company TOM Online and Skype Limited, is a Voice-over-IP (VoIP)/chat client program that uses a list of keywords to censor chat messages in either direction [35].

#### On Server System (CEN.SER).

Information can be censored on the server-side system. A censor can compromise and corrupt information on the system surreptitiously (for example by means of malware or insider attacks). A more explicit censorship policy is to make it mandatory for servers to run censorship software that ‘neutralizes’ content before it is served, or drops server responses containing blacklisted keywords. Consequently, it is possible for a censor to both corrupt information and disable its access on the server system. A number of studies investigate Chinese government’s censorship of posts on the national microblogging site Sina Weibo. Bamman et al. [15] analyze three months of Weibo data and find that 16% of politically-driven content is deleted. Zhu et al. [82] note that Weibo’s user-generated content is mainly removed during the hour following the post

<sup>4</sup>We refer to information access/distribution as access alone henceforth.

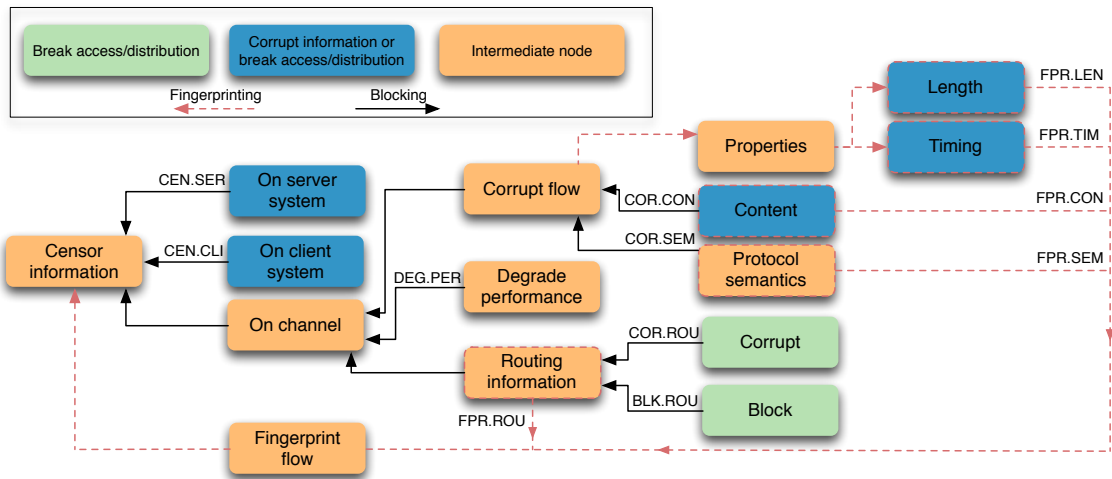


Fig. 2: Censor’s attack model, showing both direct targeting (information corruption or disabling access) and indirect targeting (fingerprinting to develop new features for information corruption or disabling access)

with  $\sim 30\%$  of removals occurring within 30 minutes and  $\sim 90\%$  within 24 hours. Another study observes posts from politically active Weibo users over 44 days and finds that censorship varies across topics, with the highest deletion rate culminating at 82%. They further note the use of *morphs*-adapted variants of words to avoid keyword-based censorship. Weiboscope [31], a data collection, image aggregation and visualization tool, makes censored Sina Weibo posts by a set of Chinese microbloggers publicly available.

#### Degrade Performance (DEG.PER).

A censor can degrade network performance as a soft form of blocking. The induced performance overhead discourages users from using a service while at the same time affording deniability to the censor. As a result, information access is hurt. Anderson [6] uses a set of diagnostics data (such as network congestion, packet loss, latency, bottlenecks) to study the use of throttling of Internet connectivity in Iran between January 2010 and 2013. He uncovers two extended periods with a 77% and 69% decrease in download throughput respectively; as well as eight to nine shorter periods. These often coincide with holidays, protest events, international political turmoil and important anniversaries, and are sometimes corroborated by overt filtering of online services or jamming of international broadcast television.

#### Block Routing Information (BLK.ROU).

While the client and server systems are not generally under control of the censor, the communication channel may be run by the government or a government-authorized telecoms provider, so blocking on the channel is generally easier than at the edges. A censor can enforce a blocking policy based on elements of the connection tuple used in routing policies: source IP address, source port, destination IP address and destination port. The block can continue for a short period of time to create a chilling effect and encourage self-censorship on part of the client. One study notes that the Great Firewall of China (GFW) has blocked communication from a client IP address to a destination IP address and port combination for

90 seconds after observing ‘objectionable’ activity over that flow [1]. It is unusual for a censor to block source port, but can potentially be used as part of a censorship policy where the source port is known to be associated with a circumvention software. To reduce collateral damage, GFW drops packets originating from Tor bridges based on both source IP address and port [76]. Blocking of this kind leads to information becoming inaccessible.

#### Corrupt Routing Information (COR.ROU).

Instead of blocking endpoints, a censor can corrupt information that supports correct routing of packets. This can be done by changing routing entries on an intermediate censor-controlled router. Alternatively, a censor can effect the same by manipulating information that supports the routing process, for example BGP hijacking and DNS manipulation. Border Gateway Protocol (BGP) is the de facto protocol for inter-AS routing. A censor can block a network’s connectivity to the Internet by withdrawing previously advertised network prefixes or re-advertising them with different properties (rogue BGP route advertisements). Many countries have attempted to effect complete or partial Internet outages in recent years by withdrawing their networks in the Internet’s global routing table (Egypt [51], Libya [52], Sudan [54], Myanmar [53]). Like BGP, DNS is another vital service over the Internet that supports its operation by mapping names given to different Internet resources to IP addresses. DNS is a hierarchical distributed system and a censor can manipulate its operation on portions of DNS that fall under its control. This may involve redirecting DNS queries for blacklisted domain names to a censor-controlled IP address (DNS redirection/poisoning), a non-existent IP address (DNS blackholing) or simply dropping DNS responses for blacklisted domains. China’s injection of forged DNS responses to queries for blocked domain names is well known, and causes large scale collateral damage by applying the same censorship policy to outside traffic that traverses Chinese links [7]. Regardless of the vector chosen to corrupt routing information, its consequence is that information access

is disabled.

#### **Corrupt Flow Content (COR.CON).**

A censor can compromise information integrity by corrupting flow content (COR.CON). In the current context, content of a flow refers to the information content as perceived by the application layer (for example an HTML document transferred over HTTP). A censor can delete, modify or insert information into original content, effectively modifying its interpretation from the one originally intended by its sender. For example, a censor can inject HTTP 404 Not Found message in response to requests for censored content and drop the original response. Alternatively, a censor can modify the body of original HTTP responses with something of its choice. Corruption of flow content implies corruption of information and/or blocking of access.

#### **Corrupt Protocol Semantics (COR.SEM).**

Another class of attacks achieves censorship by manipulating protocol semantics (COR.SEM). A censor can exploit knowledge of protocol specification to induce disruption on a flow; for example injecting forged TCP reset packets into the flow will cause both endpoints to tear down the connection. A censor can combine its knowledge of protocol semantics with manipulation of flow timing to induce flow interpretation of its choice on the destination. Consequently, corruption of protocol semantics can disrupt both information access and integrity.

### *B. Fingerprinting*

A censor needs some criteria to refine the blocking decisions described in the previous section—indiscriminately applying these to all traffic/systems would lead to the disruption of a large volume of legitimate traffic which could be unacceptable to the censor depending on various factors including cost, policy and technology [71]. Consequently, a censor performs a range of activities (*flow fingerprinting*) to scrape information (e.g. blacklist of IP addresses, keywords) to aid its blocking decision.

#### **Routing Information (FPR.ROU).**

A flow can be associated with a protocol based on elements of the connection tuple used by routing policies, namely source and destination IP addresses and ports. Destination port is a typical target of censorship (e.g. 80 for HTTP); less commonly, flows to an IP address known to be exclusively associated with a blocked service can be disrupted by implication. Flow fingerprinting of this kind can form part of a multi-stage blocking policy, possibly followed by a blocking step. Clayton examines the hybrid two-stage censorship system *CleanFeed* deployed by British ISP BT. In the first stage, it redirects suspicious traffic (based on destination IP and port) to an HTTP proxy. In the next stage it performs content filtering on the redirected traffic and returns an error message if requested content is in the Internet Watch Foundation (IWF) list [13].

#### **Content (FPR.CON).**

Another method is to inspect flows for the presence of content indicative of a protocol to be blocked, or matching a blacklist of keywords, domain names and HTTP hosts etc. A number of DPI boxes can perform regex-based traffic classification [36], [10], [60], [46], however it remains unclear what are the true

costs of performing deep packet inspection (DPI) at scale [16], [61]. Alternatively, flows can be fingerprinted based on some property of the content being carried. For example, a censor that does not allow encrypted content can block flows where content has high entropy [17].

#### **Flow Properties (FPR.LEN and FPR.TIM).**

A censor can fingerprint a protocol by creating a statistical model based on its flow features such as packet length, and timing-related features (inter-arrival times, burstiness etc). With this model, a censor can simply censor a protocol based on flow resemblance/anomaly to it [8], [78]. Wiley [74] used Bayesian models created from sample traffic to fingerprint obfuscated protocols (Dust [75], SSL, obfs-openssh [38]) based on flow features, and found that across these protocols length and timing detectors achieved accuracy of 16% and 89% respectively over entire packet streams, while the entropy detector was 94% accurate using only the first packet. Another vein of work classifies traffic flows by application using a host's transport layer behavior (such as the number of a host's outgoing connections). Flow properties can also be used to fingerprint the website a user is visiting even if the flow is encrypted [48], [63], [25], [9].

#### **Protocol Semantics(FPR.SEM).**

A censor can fingerprint flows based on protocol behaviour triggered through different kinds of active manipulation: drop, inject, modify and delay packets. Regardless of the mechanism used, the key idea is to elicit some information by leveraging knowledge of the protocol's semantic properties. If any of these techniques elicit the behaviour of a known protocol, the flow can be flagged for subsequent blocking. Alternatively, a censor can perform several fingerprinting cycles to elicit the information on which to base subsequent blocking decisions. In 2011, Wilde [64] investigated how China blocked Tor bridges and found that unpublished Tor bridges are first scanned and then blocked by GFW. Wilde's analysis showed that bridges were blocked in the following fashion: (i) When a Tor client within China connects to a Tor bridge/relay, GFW's DPI box flags the flow as potentially Tor flow, (ii) random Chinese IP addresses then connect to the bridge and try to establish a Tor connection; if it succeeds, the bridge IP/port combination is blocked.

### III. A PLUGGABLE TRANSPORT MODEL FOR LINK OBFUSCATION

We now turn our attention from a censor's attack landscape to censorship resistance systems (CRSs). The link between information client and server (*On Channel* in Figure 2 on page 3) is a frequent target of censorship and hence the assumed threat model of most CRSs. Most CRSs involve an intermediate proxy that divides the link into two parts, (i) client to proxy and (ii) proxy to server. The former is located within the censored region and is handled by CRS's *link obfuscation* module, while the latter is outside the censored region and the *proxy* module manages it. Link obfuscation, being located in the censored region, is the source of arms race between censorship and circumvention and an active area of research. We present the link obfuscation role of CRS as an abstract

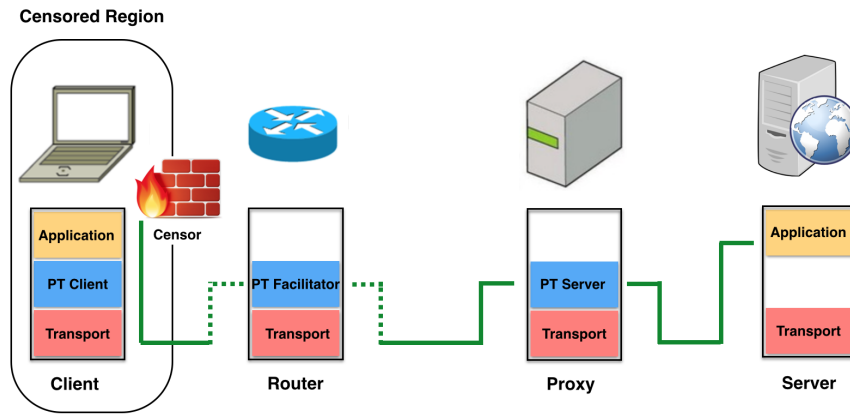


Fig. 3: A *proxy* relays traffic between a *client* in the censored region and an external *server*, and effectively divides the channel into two portions which circumvention tools handle separately via different modules: (i) client to proxy (*link obfuscation module*), and (ii) proxy to server (*proxy module*). The *Pluggable Transport (PT) client* allows the client application to construct a communication channel to the server application through a *PT server*, which can optionally be facilitated by an intermediate device (*PT Facilitator*) that partially implements the PT protocol.

model we call a *Pluggable Transport*. We further introduce an *evaluation stack* that represents functional components of a PT as a multi-layer stack. Effectively, we can use the evaluation stack as a common benchmark to visualize capabilities of different Pluggable Transports.

#### A. Background

A censor typically targets data in transit due to its less intrusive nature compared to edge-based censorship mechanisms, and because the communication infrastructure is usually under the censor’s direct control. Consequently, most censorship resistance systems (CRS) focus on circumvention of censorship on the link between information client and server. This encompasses resistance against blocking of link endpoints (IP address blocking), and fingerprinting/corruption of the content being carried over the link. To achieve these goals, CRSs typically employ a proxy, an intermediate unblocked system that relays traffic back and forth between client and server. A proxy divides the link between client and server into two distinct portions: (i) client to proxy (within censored region), and (ii) proxy to server (outside censored region). This has been illustrated in Figure 3 on page 5. A design trend is for CRSs to treat these two portions separately (via *link obfuscation module* and *proxy module*, respectively) as these lend themselves to different design, implementation, and software distribution practices.

The proxy module, being in uncensored region, may simply provide access to server, without offering any additional security properties and could be simply implemented as a HTTP or SOCKS proxy, or as a VPN. Alternatively, the proxy module may be an anonymity system like Tor [70] which not only provides access to the server but also prevents attackers from being able to identify which user is accessing which resource. Proxying is a well studied problem with widely accepted protocols for both simple proxying and anonymous communication.

In contrast to proxy module, link-obfuscation is a less mature area so it is unclear which design decisions are optimal. Efficient link obfuscation systems all have known vulnerabil-

ities which certain attackers might be capable of exploiting. There is a strong case for the public description of anonymity system designs as they can be designed in compliance with Kerckhoff’s principle – only the key must be secret in order for their security to be met. In contrast, no link-obfuscation system meets this goal while maintaining good usability, so there is a case that some link-obfuscation modules should be distributed in obfuscated binary form.

We represent the link-obfuscation role of recent CRSs as an abstract model we call a Pluggable Transport (PT). Any module which implements the abstract model can be ‘plugged-in’ to any CRS, whether it is a simple proxy or full blown anonymity system.

#### B. Abstract Model of a Pluggable Transport

The goal of a Pluggable Transport (PT) is to enable a client application to communicate with a server application over the Internet, even though direct connections are blocked. The PT-client exposes an API whereby the client application can request that a communication channel to the server application be opened. The PT-client then connects to the PT-server over a blocking-resistant communication channel, and the PT-server connects to the server application. The client application can then communicate with the PT-client, as if it is communicating directly with the server.

The communication channel provided by the PT-client and server has similar properties to TCP. Data sent through the channel will either be delivered to the other end without corruption in the same order as it was sent, or an error will be reported to the sender. It is the responsibility of the PT to route communications between the PT-client and PT-server, avoid blocking, and recover from any corruption of data (whether by the censor or due to other network disruption).

The PT-client and PT-server may be able to communicate directly over the Internet, in that any intermediate networks or routers are not aware of the protocol the PT is using. However in some cases there may be a PT-router which implements part of the PT protocol so as to facilitate the blocking resistant communication channel.

The PT communication channel does not offer authenticity, so it is the responsibility of the client and server applications to confirm that data received on the channel originated from the expected party and has not been corrupted in transit. However many practical PTs will provide some degree of authenticity so as to meet the goal of blocking resistance.

Ideally the latency of the communication channel will not be much higher than that of the direct communication channel, but in some cases a much higher latency is unavoidable. In which cases a client application which is designed for a normal TCP connection may malfunction.

This abstract model is implemented through the de-facto Pluggable Transport standard (PT) [69]. This API specifies how PT modules are invoked (i.e. as separate executable), how the communication channel is implemented (i.e. as an extension to the SOCKS protocol), how configuration and status information is communicated between client/server applications and PT client/server (i.e. through a combination of environment variables, command-line parameters, and standard in/out/error descriptors). Although this specification was written for the Tor anonymity system, it also is implemented in the Lantern [37] and Psiphon [50] simple proxy CRSs. The link-obfuscation systems discussed in the remainder of the paper comply with the abstract PT model and either implement the Pluggable Transport specification or could be adapted to do so without great difficulty.

### C. Evaluation Stack

In order to defend against the multiple avenues of attacks available to a censor, a Pluggable Transport (PT) is typically designed as a series of components, with each component defending against one or more attacks, either by itself or in conjunction with other components. In order to describe the capability of each PT, we will map each of them to a generic set of components shown in Figure 4 on page 7, arranged in layers analogous to a network protocol stack.

The primary flow of payload information is between adjacent layers in the stack, with the client/server application at the top and network at the bottom. However just as with real-world network stacks, control information does not always exactly follow this abstraction and may skip layers. Also not all layers will be present in all PTs, as certain exclude some attacks from their threat model.

An uppermost layer is **Session Initialisation (SI)**, which does not carry payload data and so is not directly connected to the application. SI is responsible for the handshake between PT-client and PT-server which may include negotiating connection parameters, performing authentication, and deriving session keys.

Also on the uppermost layer is **Encryption (ENC)** which takes application traffic which has consistent patterns and converts it to data which the adversary cannot distinguish from random. The key for performing the encryption is provided by the SI layer.

Next is the **Multiplexing (MUX)** layer, which allows multiple application channels to be multiplexed over a single PT channel, or a single application channel to be split over

multiple PT channels. This layer is also responsible for error detection and re-assembly, if the lower layers do not provide this.

Then **Content Obfuscation (OBF)** transforms the ‘random’ data into traffic which appears to be a different protocol, e.g. HTTP or VoIP.

Next **Timing Obfuscation and Length Obfuscation (TIM-LEN)** hide the application’s timing and packet length patterns. The layer may perform the obfuscation itself, or may just compute the changes which are necessary and transmit this as control data to other layers which actually delay and/or pad payload data.

Finally **Transport (TRN)** is responsible for taking the transformed payload data and sending it to the other side of the PT client/server pair.

We survey 34 existing CRSs that fit our abstract PT model and broadly systematise them according to the threat landscape these assume, i.e. the path(s) on the attack model (Figure 2 on page 3) that they protect. We describe six Pluggable Transports at length using our evaluation stack. We believe that the PTs we select are a good representation of their corresponding class. We review multiple PTs within a class to reflect the diversity of implementation, and provide brief description of other tools to show breadth of the area.

## IV. IP ADDRESS/HOST FILTERING RESISTANCE SYSTEMS

A censor can disrupt access to a service by blocking the IP address of the server, or a key hop directly on path to the server (corresponding to BLK.ROU in Figure 2 on page 3). A number of tools have emerged to resist IP address filtering, of which proxies are the most prevalent. A proxy relays traffic between the source and destination effectively obfuscating the latter [72], [19]. However, simple proxies with long-lived IP addresses can be easily enumerated and blocked by a censor, thus motivating development of more sophisticated mechanisms to resist IP address, domain name and host filtering. We classify existing work into two categories based on the approach they take to achieve circumvention. Techniques under *sensorship surface augmentation* hide the censorship target (e.g. IP address) among a crowd such that censoring the target incurs larger collateral damage than if it is blocked in isolation. In *decoy routing*, a cooperative hop between client and server applies circumvention-friendly treatment to packets containing a special steganographic mark.

### A. Sensorship Surface Augmentation

A censor’s blocking decision has associated accuracy depending on a number of factors, such as the blocking mechanism employed and the quality of target set to block (e.g. how representative is a blacklist of IP addresses of a block category such as ‘porn’). In particular, the censor’s policy must make consideration for the acceptable false positive rate as these have political and economic ramifications. Mechanisms under this category leverage this observation by obfuscating the censorship target in such a way that censoring it incurs large collateral damage in terms of false positives.

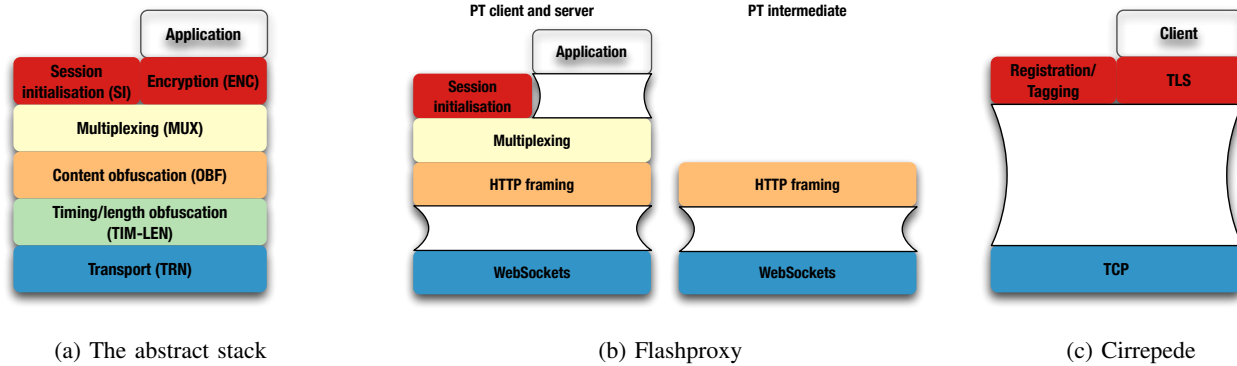


Fig. 4: The Abstract Stack and Evaluation Stacks for Flashproxy and Cirrepede.

1) *Flashproxy*: A popular means to resist IP address filtering is to proxy traffic between client and blocked server through an intermediate host, this effectively hides the IP address and name of a blocked server. By pretending to be a genuine client, a censor could enumerate long-lived proxies and subsequently block them (FPR.ROU, COR.ROU and BLK.ROU nodes in the attack diagram in Figure 2 on page 3).

To protect against these attacks, Flashproxy introduces two new entities, i.e. *facilitators* (Figure 4b on page 7, left) and *flashproxies* (Figure 4b on page 7, right). A *facilitator* is a volunteer website outside the filtered region which may be blocked by the censor, yet remains reachable through low-bandwidth channels such as email. Over this channel, the *Session Initialisation* module (*SI*) of a censored user registers itself; that is, it sends its IP address and a port where it awaits incoming connections. On the *facilitator* side, the *SI* adds a special *badge* in all the web pages it serves to uncensored users, typically a piece of javascript. When an uncensored user visits a web page on the *facilitator* website, the *badge* turns his web browser into a so-called *flashproxy*. The *SI* of a *flashproxy* (running in the web browser of an uncensored user) connects to the *facilitator* to retrieve an IP/port pair for a censored user and initiates a connection to it. The *SI* of the censored user accepts the connection to complete the rendezvous between the two entities. Thereafter, the *Content Obfuscation* (*OBF*) of the *flashproxy* relays censored content for the censored user through HTTP (“HTTP framing” in Figure 4b on page 7).

Flashproxy offers resistance against IP address blocking only, consequently leaving a number of paths on the attack diagram exposed. A censor could block traffic based on content (COR.CON), use statistical traffic properties (FPR.LEN and FPR.TIM) to detect the protocol or content to block, or observe characteristic patterns in incoming connections to censored hosts (FPR.SEM). The authors suggest using Flashproxy in combination with Tor to thwart COR.CON attacks.

2) *Others*: One design trend is to use a widely used service as a proxy to fetch censored content. Meek [43] employs a technique called *domain fronting* to evade host based censorship by using an innocuous domain name (*front domain*) in the unencrypted request header (*TLS Server Name Indication header* – SNI), while hiding the domain of a proxy

(*inside-domain*) in the encapsulated encrypted request (*HTTP Host* header). The front domain is an intermediate web service hosting many domains (typically a CDN) which decrypts the inside-domain and internally routes the traffic to the relevant host within its network. This host serves as a proxy for censored clients to access blocked servers. OSS [21] turns any existing Online Scanning Service (OSS) into a proxy. These are web services that take a URL as user input and then fetch the web page behind that URL (e.g. PDFmyURL [49]). A censored client Alice transmits a request to a non-blocked machine Bob by providing the OSS with a URL such as *www.bob.com/censored-request*. This makes the OSS connect to Bob with the *censored-request* encoded in the URL. Bob responds using a redirection mechanism such as *HTTP 302 Found* with a *Location* header of the form *www.alice.com/censored-response* that points back to Alice. If the OSS follows redirections, it connects back to Alice with the *censored-response* embedded in the url. A variation of this scheme is for clients and servers to rendezvous on an intermediate host, blocking which incurs significant collateral damage. CloudTransport [44] clients and bridges share a common account on a cloud storage which they use to share files containing client requests and bridge responses in real time. Collage [11] peers exchange data through social networking and photo sharing websites by embedding hidden messages into user-generated content such as posts and images. The popularity of these websites and the bulk of volume generated on them makes it a hard for a censor to accurately spot and block censored content. MIAB [30] improves Collage’s rendezvous by leveraging *blog pings*, i.e. real-time notifications a blog sends to a centralized network service (a ping server) when content is updated. By monitoring ping servers, MIAB peers automatically learn within minutes when a new message is available. Defiance [40] allocates ephemeral IP addresses to its gateways and bridges from a large pool of diverse IP addresses. The transiency and diversity of IP addresses make it hard for a censor to block or enumerate them. To access blocked websites, a Defiance client must connect to a shorted-lived bridge which acts as a proxy. To learn an ephemeral bridge location, a client must successfully complete a *dance*; that is, make a sequence of pre-agreed timed short-lived connections to ephemeral gateways.

## B. Decoy Routing

This approach resists IP address filtering by having clients covertly signal a cooperating router along the way to deflect their traffic intended to a non-blocked destination to a blocked one. This thwarts attacks on the FPR.ROU, COR.ROU and BLK.ROU nodes. To deflect traffic, deflecting routers must be located on the forward network path from the client to the non-blocked destination. These routers must therefore be strategically positioned to optimise the number of censored users that can be served [12]. It is theoretically possible for a censor to defeat decoy routing by routing traffic around the deflecting routers [57]. However, in practice this is believed to be too costly for a censor because of business relationships with other ISPs, and monetary, performance and quality of service degradation issues thereby induced [28].

1) *Cirrepede*: Its main contribution is resistance against IP address filtering (BLK.ROU, FPR.ROU and COR.ROU). In addition to this, it also resists content-based fingerprinting (FPR.CON) and tampering/blocking (COR.CON) through its authenticated encryption layer. Its corresponding evaluation stack is presented in Figure 4c on page 7.

To use Cirrepede, a censored user must first register its IP address and a shared secret with a Registration Server (RS). A client never “talks” directly to this server. Instead, a friendly ISP deploys Deflecting Routers (DRs) that redirect non-registered client traffic to the Registration Server. To register, the *Session Initialisation (SI)*, *Transport (TRN)* and *Encryption (ENC)* module of a client coordinate to encode a covert registration signal into the *TCP Initial Sequence Number (ISN)* of a series of packets destined to a non-blocked destination. The registration packets are deflected by the Deflecting Routers (DRs) and reach the Registration Server where they are inspected. If the Registration Server successfully recognises the signal in the packets, its *SI* instructs all Deflecting Routers within the ISP network to deflect subsequent client traffic to a Service Proxy (SP). The registration is opportunistic in the sense that there is a chance that no Deflecting Router is located on the forward path between the client and the non-blocked destination; and therefore that the registration packets never reach the Registration Server.

After sending the registration packets, a client selects an innocuous non-blocked destination and initiates a TLS handshake to it. The packets are deflected by the Deflecting Routers to the Service Proxy which observes the communication. When the TLS handshake completes, the Service Proxy takes over the connection; that is, it closes the TCP connection with the decoy server on behalf of the client whilst keeping the client side of the connection open. Importantly, it blocks all subsequent traffic from the decoy server to the client so as not to arouse the censor’s suspicion. This is sometimes referred to as “inline flow blocking” because the Service Proxy must be between the client and the decoy destination. At this point, the *SI* of the client and Service Proxy derive a new cryptographic key based on the shared secret exchanged during registration. To ensure both parties have successfully derived the new key, the *SI* of the Service Proxy encrypts a known value and sends it to the client over the existing TCP connection. If the

client successfully decrypts it, the new key is kept and used to encrypt further communication between the two parties. Thereafter, the Service Proxy acts as a web proxy to censored content.

A censor could still attack unprotected nodes of the attack diagram. Traffic from/to different websites generally have different characteristic patterns, therefore a censor could determine that traffic allegedly originating from a non-blocked website comes from a different (FPR.LEN and FPR.TIM). It could also detect protocol implementation inconsistencies due to the non-blocked and blocked destination running different software stacks (FPR.SEM). Cirrepede is also vulnerable to replay attacks: a censor who replays the registration packets will find itself incapable of decrypting the TLS traffic after the handshake: this hints that a new key is covertly negotiated during the TLS handshake.

2) *Others*: Telex [80], TapDance [79] and Curveball [32], [33] can selectively *tag* individual connections on-the-fly (in contrast to Cirrepede that deflects client traffic *after* registration). Telex and Curveball embed their tag in the random nonce of a TLS handshake with a non-blocked destination. TapDance encodes it in a connection’s incomplete HTTPS request destined to a decoy server using a novel steganography scheme. As Cirrepede, TapDance and Curveball support asymmetric flows where upstream and downstream traffic do not follow the same path because of ISP’s internal routing. TapDance is the only solution that does not require active “inline flow blocking” to prevent further decoy-server client communication, for this reason it is believed to be easier to deploy by ISPs without disturbing existing traffic and quality of service. IBS [55] improves decoy routing solutions in general by simplifying key distribution and providing forward secrecy. These are achieved with the use of identity-based encryption instead of traditional public-key cryptography.

## V. FLOW FINGERPRINTING RESISTANCE SYSTEMS

There has been significant work on obfuscating blocked traffic such that it evades a censor’s protocol fingerprinting machinery. These approaches can be divided into ones that provide resistance against *fingerprinting of protocol semantics*, while others *mimic* a supposedly whitelisted category. Some tools transform a flow such that the censor’s analysis limitations fail to trigger censorship (*monitor-driven flow transformation*).

### A. Fingerprinting of Protocol Semantics

A number of schemes offers resistance against protocol scanning. These are techniques currently used by censors to confirm that a server is indeed part of an anti-censorship system. Typically, a censor would probe for an open port and attempt to “speak” the anti-censorship protocol. To defeat protocol scanning, a SilentKnock [73] server accepts incoming connections on a particular port only from clients that authenticate with a special “knock”. This knock is a one-way authentication mechanism embedded in TCP headers, it is indistinguishable from an ordinary TCP/IP connection and



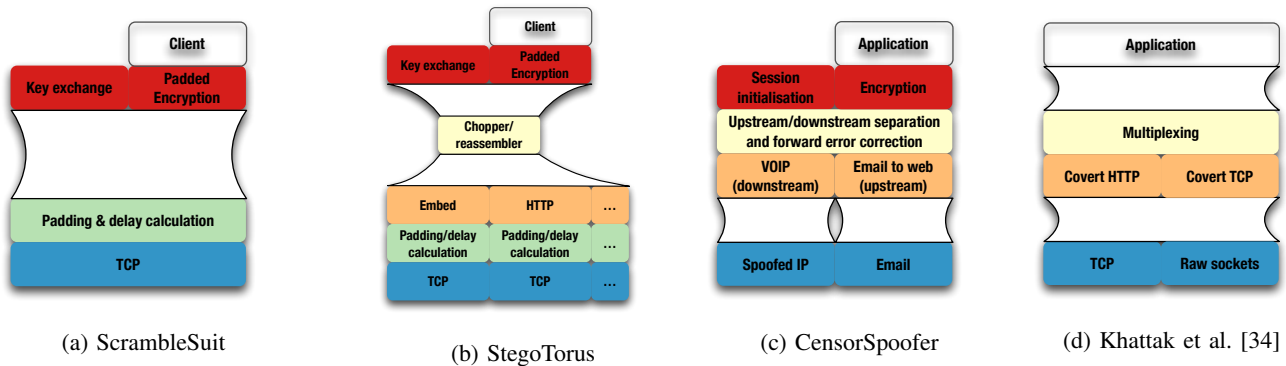


Fig. 5: Evaluation Stacks for ScrambleSuit, StegoTorus, CensorSpoofer and Khattak et al. [34].

resists forgery and replay attacks. BridgeSPA [59] (originally known as SPATor [58]) builds upon SilentKnock’s design but relaxes server-side memory constraints associated with per-client housekeeping. It replaces counters with rounded-to-the-minute UTC timestamps and long-lived keys with short-lived ones. Defiance [40] imposes several levels of address indirection to prevent unauthenticated access to bridges from subsequent protocol fingerprinting.

## B. Mimicry

Mimicry-based mechanisms transform traffic to look like whitelisted communication. The transformation can be applied at both transport and application layers, and the transformed traffic resembles the syntax and/or content of an allowed protocol, or randomness. Additionally, some of these solutions also resist protocol fingerprinting and scanning.

1) *Mimic Existing Protocol or Content:* A large body of work evades censorship by imitating an innocuous protocol (e.g. HTTP) or content (e.g. HTML). By mimicking widely deployed protocols and content, a PT increases the censor’s work load as there is more volume of traffic to inspect. Furthermore, a censor is reluctant to conduct wholesale censorship of a popular protocol/content due to the associated collateral damage.

### StegoTorus.

StegoTorus obfuscates Tor packet lengths (FPR.LEN) and inter-packet timings (FPR.TIM). It also prevents content fingerprinting (FPR.CON) and tampering/blocking (COR.CON) with authenticated encryption. Optionally it can mimic a set of innocuous protocols over which covert traffic is sent. The corresponding evaluation stack is presented in Figure 5b on page 9.

StegoTorus comprises two modules, both of which can be implemented by a combination of PT components. To start a session, the *SI* of a client and server (i.e. proxy) first establish a shared key through a key-exchange that only contains random bytes, this thwarts content-based filtering (FPR.CON). Once a session is established, the *TIM-LEN* module chops fixed-length input packets into random-sized messages. Sizes are taken from a trace of an innocuous protocol session pre-recorded by the user or bundled with the software (FPR.LEN). A 32-byte ID is added to messages so they can be re-ordered

by the recipient. The *ENC* module then encrypts messages individually (COR.CON and FPR.CON) and passes them again to the *TIM-LEN* module which adjusts their sending time in accordance with packet inter-arrival times derived from a pre-recorded traffic trace (FPR.TIM). Optionally, the *Content Obfuscation (OBF)* component of StegoTorus can mimic a known protocol such as unencrypted HTTP. For example the *OBF* of a client may embed the encrypted messages into the HTTP *Cookie* header and url part of requests, and the *OBF* of a server may reply by steganographically embedding content into the body of HTTP responses.

StegoTorus uses long-lived server locations. So if a censor manages to harvest them, it could block them by injecting TCP *RST* packets or erroneous DNS responses (COR.ROU, BLK.ROU and FPR.ROU).

### Others.

FOE [22] and MailMyWeb [41] proxies send users static web pages as email attachments in response to censored URLs they receive in an email *Subject* field. These solutions assume that SMTP is not monitored or email traffic is encrypted. SWEET [81] creates a bi-directional communication channel by encapsulating censored traffic into email attachments such as images. To prevent a censor from blocking all emails sent to the proxy, each client sends requests to a unique email address. SkypeMorph [45] shapes inter-packet timing and packet size distribution so as to mimic a Skype video call. The shaping resists both low and high order statistical inference attacks. TransTeg [42] negotiates an *overt codec* during a VoIP call initialisation, but thereafter encodes the raw voice stream with a different lower-bitrate codec (the *covert codec*). This effectively reduces the length of packets to transmit. The space freed is filled with low-bandwidth covert traffic. FTE [18] extends conventional symmetric encryption with the ability to specify the format of the ciphertext with a regex. This effectively transforms a blocked *source* application-layer protocol into an unblocked *target* application-layer protocol, so that DPI boxes mis-identify a blocked protocol as a non-blocked one.

Protocol imitation has a number of limitations that make it possible to distinguish imitated traffic from legitimate traffic. An ideal protocol imitation must not only adhere to its specification, but also mimick (*i*) other protocols it depends on (e.g. HTTP relies on DNS to find a host IP address),

(ii) dependencies between multiple connections (e.g. a VoIP must start with SIP messages and be followed by UDP traffic and RTP/RTCP connections in that order), (iii) reaction to network disruption (e.g. packet congestion), (iv) geolocation optimisations (e.g. Web services minimise round-trip time by locating their servers close to clients), and (v) software artefacts introduced by specific implementation stacks, operating systems or programs. Houmansadr et al. [26] therefore conclude that protocol imitation is a flawed approach because a comprehensive mimicry must exhibit every observable aspect of it, but this is challenging, time-consuming and error prone. Geddes et al. [23] show that solutions that embed covert traffic into voice and video streams are inherently flawed because of fundamental differences in channel requirements of the overt and covert protocol. The overt protocol is peer-to-peer and loss tolerant while the covert one is client-proxy and loss intolerant. Consequently, a censor can disrupt the covert communication without affecting legitimate traffic. To avoid the above pitfalls, another approach to mimicry is to re-use genuine software and libraries and tunnel covert traffic through them. Facet [39] streams video over a video-conferencing calls such as Skype, Google Hangout or FaceTime calls. Client requests are sent to a Facet server through a low-bandwidth upstream channel such as email or IM. In response, the server initiates a voice-conferencing session and streams the censored video. Freewave [27] modulates traffic into acoustic signals and streams them directly into an existing VoIP application such as Skype.

2) *Mimic Unknown Protocol or Content*: Another approach is to make traffic look like an unknown protocol, either by imitating randomness or arbitrarily deviating from a known blocked one. This idea is motivated by the general assumption that a censor implements blacklisting of known protocols and is unwilling to incur high collateral damage associated with whitelisting.

#### **ScrambleSuit.**

ScrambleSuit [77] obfuscates a censored protocol with random-looking bytes for all its traffic including during session initialisation (FPR.CON). The encryption provides content obfuscation (FPR.CON), confidentiality and resistance against tampering (COR.CON). Packet lengths and inter-arrival times are also randomised (FPR.LEN and FPR.TIM). The protocol used to bootstrap a session resists tampering (COR.CON) and protocol scanning (FPR.SEM). Its corresponding evaluation stack is presented in Figure 5a on page 9.

To connect to a ScrambleSuit proxy, the *SI* of a client redeems a short-lived ticket retrieved from a low-bandwidth out-of-band channel. After the first authentication, the server gives the client a ticket for the next connection; so it need not retrieve one out-of-band again. This ticket provides mutual authentication and therefore resistance against protocol scanning from a censor (FPR.SEM). Furthermore, it only contains random bytes to evade content-based detection (FPR.CON). Once a session is initialised, the *ENC* component turns all traffic into random-looking bytes, thereby thwarting content-based fingerprinting and blocking (FPR.CON). It also provides authentication and confidentiality, except in systems where multiple clients share the same bridge (COR.CON). The

*TIM-LEN* component randomises packet lengths (FPR.LEN) and timing of flows (FPR.TIM) using discrete probability distributions provided by the *ENC* module.

Despite protecting many nodes in the attack diagram, ScrambleSuit still leaves some of them for the censor to attack on. For example, it does not hide the IP addresses of proxies. Therefore, if a censor manages to learn them, it could subsequently block them (COR.ROU and BLK.ROU).

#### **Others.**

MSE [56] is the de-facto obfuscation mechanism used by BitTorrent. The key exchange and traffic contain only random-looking bytes. Padding is added to all traffic including the handshake. Dust [75] follows the same methodology but can shape packet sizes based on an arbitrary distribution. obfs2 [67] also obfuscates traffic content with encryption, but the key exchange does not provide authentication against passive and active attackers. obfs3 [68], initially adopted by ScrambleSuit, improves obfs2 by negotiating keys using anonymous Diffie Hellman (DH) with a special encoding so as to be indistinguishable from a random string. This forces a censor to actively probe the server or perform a man-in-the-middle to detect the key exchange. obfs4 [66] adds authentication to obfs3 using the ntor handshake [24], [47] and is now used by ScrambleSuit. Unlike ScrambleSuit, obfs2 and obfs3 do not randomise packet lengths and inter-packet timings.

### *C. Monitor-driven Flow Transformation*

Censorship systems have an analysis model to identify targets to filter (flow fingerprinting). The model employed likely has some limitation with respect to traffic analysis. Tools under this category shape traffic such that it evades a censor due to limitations in its traffic analysis model. Effectively, this mechanism can provide unobservability to even cleartext traffic which is particularly useful for countries that disallow encrypted traffic.

1) *Khattak et al. [34]*: This work explores protection against content-based filtering (FPR.CON and COR.CON) by exploiting characteristics of the implementation of DPI boxes. The evaluation stack is presented in Figure 5d on page 9.

Since DPI boxes operate in the same manner as Network Intrusion Detection Systems (NIDS), they are vulnerable to the same evasion techniques. The authors treat DPI boxes as black-boxes. They assume a set of evasion classes that DPI boxes may be vulnerable to, and test their hypotheses by sending specially crafted packets. A DPI box typically operates in the following order to keep track of traffic. First it identifies new TCP flows based on a full or partial TCP handshake: this triggers the creation of a Transmission Control Block (TCB) in memory to keep a flow's current state and properties. Subsequent packets belonging to this stream can be examined for content. A box discards its TCB upon TCP tear-down, typically in the presence of an *RST* packet or *FIN-ACK* exchange. However there may also be connectivity issues that can lead to a connection tear-down which the box must account for. Furthermore, a box generally has an incomplete view of a flow since packets within the same stream may be

routed through different paths and therefore be unobservable. Other inaccuracies can arise if the implementation does not validate header fields or when messages diverging from the protocol specification are observed. Overlapping IP fragments and TCP segments must also be re-constructed under certain CPU and latency constraints. DPI boxes from different vendors take different views on these, but it is generally a source of vulnerabilities. As a proof-of-concept, the authors perform a series of tests on GFW. They find that by sending an *RST* packet with a low *TTL* value for an existing connection, subsequent packets containing censored keywords are received by the end point. This confirms that the DPI box has freed its TCB for the flow even if the *RST* packet is never received by the destination.

Most nodes of the attack diagram are not protected by this approach. Even without a clear view of the entire stream, a censor could still block traffic on a per-packet basis, for example by filtering those destined to known server IP addresses (FPR.ROU, COR.ROU and BLK.ROU) or those with a particular length such as Tor’s (FPR.LEN).

2) *Others*: Clayton et al.[14] note that GFW terminates a connection containing blacklisted keywords by sending to both ends packets with TCP *RST* flag. A client can circumvent this by ignoring all *RST* packets received. West Chamber [2], [3], [4] evades GFW by misleading it to believe that a connection has terminated. To achieve this, end points exchange specially crafted *RST* packets that are ignored by their TCP stack but considered by GFW.

## VI. COMPOSITE CENSORSHIP RESISTANCE SYSTEMS

These systems offer resistance against filtering of both IP address/host (Section IV) and flows (Section V).

### A. *CensorSpoofers*

CensorSpoofers provides resistance against IP address harvesting (BLK.ROU, COR.ROU, FPR.ROU) using spoofed source IP address, and resistance against content-based blocking (FPR.CON and COR.CON) by mimicking encrypted VoIP traffic. It does not reveal the IP address of proxies to clients. The evaluation stack is presented in Figure 5c on page 9.

To use the system, a client must first register with the CensorSpoofers proxy (called the *spoofers*). For this, the *SI* of a client creates four accounts: two email accounts (one from a local provider and one from a foreign one), and two VoIP accounts (one from a local registrar and one from a foreign one). It creates a registration message containing the credentials of the two foreign accounts, the addresses of the local accounts (without credentials), and a shared cryptographic key. To complete the registration, the client must ask an existing client to send the registration on his behalf to the *spoofers*. Upon reception of the registration, the *spoofers* logs in the foreign accounts (VoIP and email) and starts monitoring incoming messages from the local accounts (i.e. from the client).

Once registered, the *SI* of a client initiates a session by calling the foreign VoIP account (which is used by the

*spoofers*). The *spoofers* accepts the call and provides its location (IP address and port) where it awaits incoming UDP voice packets. The subterfuge is that it does not provide its real IP address but a dummy one that belongs to a device on the Internet. Upon reception of the reply, the client starts sending decoy UDP traffic to the IP address provided by the censor, as would a typical client do. This traffic is never received by the *spoofers* though: it reaches the dummy host who simply ignores it without replying. This works because Internet standards require a host not to send reply packets to incoming packets on a port that is not “closed”. Meanwhile, the client sends censored requests via email (to the foreign email address registered). The *spoofers* receives the email and sends the encrypted response in UDP traffic by spoofing the source IP address of the dummy host. The session is thereby established: the censor thinks the censored user is in a voice call with an innocuous person outside the filtered region; whereas it is covertly communicating with the *spoofers*. The *ENC* module of CensorSpoofers prevents content-based detection (FPR.CON) and provides data confidentiality and authentication (COR.CON).

CensorSpoofers does not shape traffic, therefore a censor may be able to detect discrepancies between a real voice call and the web traffic sent over UDP (FPR.LEN and FPR.TIM). Furthermore, there may be an indicative correlation between the time of a VoIP call and the time of STMP traffic (FPR.SEM).

### B. *Others*

Freewave [27] mimics VoIP traffic whilst hiding proxy IP addresses. It modulates traffic into acoustic signals and sends them over a VoIP network such as Skype, Vonage or iCal. To hide the IP address of a proxy, client traffic is relayed through multiple VoIP peers. In the case of Skype, this is achieved by configuring the Freewave proxy as an ordinary node so that VoIP traffic is routed via Skype *super nodes*. Infranet [20] thwarts IP blocking by turning a non-blocked cooperative website into a proxy. Besides relaying censored content, proxy websites continue to serve their usual uncensored content. For its upstream channel, Infranet encodes covert traffic in sequences of HTTP requests; for its downstream channel, it steganographically embeds content in uncensored images.

## VII. DISCUSSION AND CHALLENGES

Effective circumvention should offer resistance against all paths on the attack diagram (Figure 2 on page 3) for a given component in the information system. In particular we focus on channel-based systems as we note in Section II that this is the primary focus of censorship and by implication that of most circumvention tools. However, a concise summary of all PTs (Table I on page 12) reveals that this is typically not the case and there is a tendency for tools to cluster around resistance against either IP address/host filtering or flow fingerprinting. We also note that no current PT provides protection against the manipulation of protocol semantics (COR.SEM), maybe because it is hard and because we do not know how to model this type of resistance correctly. We

TABLE I: Summary of Pluggable Transports. Columns represent a node of the attack diagram (Figure 2 on page 3) and a symbol in a column means that the tool protects this node. To achieve the protection, a tool uses a combination of components which we represent with a different symbol as follows: ▲ is the *Session Initialisation* (SI), ♠ is the *Encryption* (ENC), ⚡ is the *Multiplexing* layer (MUX), ■ is the *Content Obfuscation* (OBF), ◆ is the *Timing Obfuscation*, ▼ is the *Length Obfuscation* and ★ is the *Transport* layer.

Section		Blocking				Fingerprinting				
		COR.CON	COR.SEM	COR.ROU	BLK.ROU	FPR.LEN	FPR.TIM	FPR.SEM	FPR.ROU	FPR.CON
Infranet	Composite VI	▲	▲		★					★
MSE	Mimicry V-B				★					
FOE	Mimicry V				★					
MailMyWeb	Mimicry V	▲								★
SilentKnock	Protocol									★
	Fingerprinting V-A							▲		★
Traffic Morphing	Mimicry V-B									
Collage	Composite VI	▲	▲		★					★
Curpede	IP Filtering IV	▲	▲	▲	▲				▲	▲
Curveball	IP Filtering IV	▲	▲	▲	★				▲	▲
TransTag	Mimicry V	▲	▲							★
SPATor	Protocol									
	Fingerprinting V-A							▲		★
BridgeSPA	Protocol									
	Fingerprinting V-A							▲		★
Telex	IP Filtering IV	▲	▲	▲	★				▲	▲
Dust	Mimicry V-B				▲				▲	▲
Defiance	IP Filtering IV	▲	▲	▲	▲				▲	▲
SkyeMorph	Mimicry V	▲	▲	■					■	
StegoTorus	Mimicry V	▲	▲							■
obfs2	Mimicry V	▲	▲							■
CensorSpoofer	Composite VI	▲	▲		★					★
Flashproxy	IP Filtering IV	▲	▲	▲	★				▲	▲
OSS	IP Filtering IV	▲	▲	▲	★				▲	▲
Firewave	Composite VI	▲	▲	▲	★				▲	▲
MIAB	IP Filtering IV	▲	▲	▲	★				▲	▲
IRS	IP Filtering IV	▲	▲	▲	★				▲	▲
ScrambleSuit	Composite VI	▲	▲	▲					▲	▲
SWEET	Mimicry V	▲	▲		★					★
FTE	Mimicry V	▲	▲							
obfs3	Mimicry V	▲	▲							
Khattak	Monitor-driven V-C							▲		■
TRIST	Mimicry V	▲	▲							■
Meek	IP Filtering IV	▲	▲	▲	★				▲	▲
TapDance	IP Filtering IV	▲	▲	▲	★				▲	▲
CloudTransport	IP Filtering IV	▲	▲	▲	★				▲	▲
obfs4	Composite VI	▲	▲	▲					▲	▲

also note that due to their monolithic design, these systems do not lend themselves very well to sharing and modularization.

Recently, there has been a case for combining multiple Pluggable Transports with the goal to increase the censorship/attack paths than those covered by either in isolation. Fog [65] uses multiple proxies to chain together Pluggable Transports in a black box fashion. This approach is not suitable for practical deployment due to a number of limitations, for example, not all combinations of Pluggable Transports make sense: the chain obfs3 (flow fingerprinting resistance) followed by Flashproxy (IP address filtering resistance) will offer more comprehensive resistance. The reverse, that is Flashproxy followed by obfs3 breaks the former’s network layer assumptions. Consequently, sharing of Pluggable Transport features has happened not in a black-box way, but through the sharing of source code. For example, LibFTE is in use by Tor (in its fteproxy Pluggable Transport form) and a number of other projects. Similarly, Meek which was originally developed for Tor now also exists in a fork by Psiphon [50] with minor adaptations.

The evaluation stack described in Section III-C provides a more systematic way to develop Pluggable Transports. Components from a Pluggable Transport can be extracted so that each component complies with the abstract model the stack defines. This approach assists the design process by providing a set of patterns to follow, and a methodology for evaluating the censorship resistance features which are offered. Just as abstractions for components have been developed for compiler design (lexer, parser, code generator) or GUI design (model, view, controller) a systematic approach to design reduces development time and improves quality of code.

Speed of development is particularly important for censorship resistance because there is no one approach which is optimally efficient and resistant to all attackers. Therefore different systems are necessary for different situations, and if the situation changes a new system may become necessary. The ability to quickly develop censorship resistance systems for particular locations (spatial agility) and in response to changes

(temporal ability) will increase the number of users of systems and the time period that their access will be interrupted. To achieve this we propose an extension of Pluggable Transports – *Tweakable Transports*.

Tweakable Transports are Pluggable Transports built out of re-usable components following the evaluation stack architecture. Each component can be replaced with another which is compatible and components can be inserted or removed. This approach allows code-reuse because a component developed for one Tweakable Transport can be used for another. In so doing, more collaboration opportunities are allowed, better testing can be performed on frequently required components improving reliability and both spatial and temporal agility.

As a result Tweakable Transports exponentially increases the number of possible link obfuscation scheme. The development effort to add one component adds not just one new Pluggable Transport, but creates a whole new family of schemes, each one of which the censor will need to test against any proposed fingerprinting or blocking technique. The increased development cost and possibility of false-positives will reduce the likelihood that a censor will be able to effectively block the resulting link-obfuscation schemes.

A particular instance of a Tweakable stack may be designed by an expert familiar with the properties of each component and a censor’s blocking techniques and so allow the trade-off between performance and censorship resistance. Alternatively instances could be automatically generated and tested against the real censorship system or a simulation of one, so as to quickly find an adequate link-obfuscation scheme.

From the Pluggable Transports summarised in this paper, adapting their implementations to be Tweakable Transports allows weaknesses to be addressed. Missing layers (e.g. resistance to timing and packet-length fingerprinting) leave some schemes open to attack. Rather than developing a component from scratch, a component can be imported from another Pluggable Transport. New schemes can also be created, such as combining the content obfuscation and timing/length ob-

fuscation with a common session initialisation and encryption component, via a multiplexing component.

In an ideal layered model, communication would only be between adjacent layers and a component need only be aware of the abstract API that each layer makes available. Practical systems rarely meet this goal, and so while the transmission of content will only be between adjacent layers, there will be control information which might violate the layer boundaries. A similar approach is taken for other systems, as for example socket options for the UNIX socket model, explicit congestion notification in network protocols, and annotations in compilers.

A challenge in designing an API for building Tweakable Transports is how to constrain the model to allow code-reuse but to still permit cross-layer communication when necessary (e.g. for traffic shaping). The approach taken by the Tweakable Transport architecture specification is for each component to define one or more upward endpoints and zero or more downward endpoints, for connection to the component directly above and directly below in the stack. Components are then arranged in a directed acyclic graph by connecting these endpoints, forming the data plane.

In addition to the data plane, the control plane is a broadcast channel by which any component can signal a JSON formatted event to any (or all) other components. Control events may be local to the host at which it was generated (e.g. for controlling the timing characteristics of data which is to be emitted) or must be transmitted to the other host (e.g. for defining a protocol which is to be impersonated). In the latter case a multiplexing layer must be defined that receives control traffic, encodes it, and transmits it over the data plane endpoint.

To support the above model, while allowing the development of innovative components, control and data messages are sent between components over a single event queue on each host. A single event queue ensures that the order of control and data events are ordered. Components are configured on initialisation with the component graph, and tag events with a source and destination, with the destination component consuming the event. However all components can see all events to allow them to monitor the performance and behaviour of other components.

Further details can be found in the Tweakable Transport API specification [62].

## VIII. CONCLUSIONS

This paper has presented a model of both attacker and defender for censorship resistance schemes. As common censorship techniques have focussed on disrupting the channel between client and server, censorship resistance schemes have been built around the *pluggable transport* abstract model. This model has proved useful in building up a diverse ecosystem of link-obfuscation schemes, but there is no one solution which offers the optimal choice for all censorship scenarios. The paper has surveyed the field of Pluggable Transports, in terms of the threat-model defended against and their evaluation in terms of an abstract protocol stack. This evaluation has led to a new design for Pluggable Transports – the Tweakable

Transport: a tool for efficiently building and evaluating a wide range of Pluggable Transports so as to increase the difficulty and cost of reliably censoring the communication channel.

## ACKNOWLEDGEMENTS

The authors thank David Fifield, Vern Paxson and Philipp Winter for their valuable suggestions and comments. This work was supported by the Engineering and Physical Sciences Research Council [grant number EP/L003406/1]; and the Royal Society [grant number UF110392].

## REFERENCES

- [1] HTTP URL/keyword detection in depth. <http://gfwrev.blogspot.jp/2010/03/http-url.html>, 2010.
- [2] Scholar Zhang: Intrusion detection evasion and black box mechanism research of the Great Firewall of China. <https://code.google.com/p/scholarzhang/>, 2010.
- [3] west-chamber-season-2. <https://code.google.com/p/west-chamber-season-2/>, 2010.
- [4] west-chamber-season-3. <https://github.com/liruqi/west-chamber-season-3/>, 2011.
- [5] AdversaryLab. . Online, November 2014. <http://www.adversarylab.org/>.
- [6] C. Anderson. Dimming the Internet: Detecting Throttling as a Mechanism of Censorship in Iran. Technical report, University of Pennsylvania, 2013.
- [7] Anonymous. The Collateral Damage of Internet Censorship by DNS Injection. *SIGCOMM Comput. Commun. Rev.*, 42(3):21–27, June 2012.
- [8] L. Bernaille and R. Teixeira. Early recognition of encrypted applications. In *Proceedings of the 8th International Conference on Passive and Active Network Measurement, PAM'07*, pages 165–175. Springer-Verlag, 2007.
- [9] G. D. Bissias, M. Liberatore, and B. N. Levine. Privacy vulnerabilities in encrypted http streams. In *Proceedings of Privacy Enhancing Technologies workshop (PET 2005)*, pages 1–11, May 2005.
- [10] Bro. <https://www.bro.org>. Online. November, 2014.
- [11] S. Burnett, N. Feamster, and S. Vempala. Chipping Away at Censorship Firewalls with User-Generated Content. In *USENIX Security Symposium*, Washington, DC, USA, 2010. USENIX.
- [12] J. Cesareo, J. Karlin, M. Schapira, and J. Rexford. Optimizing the placement of implicit proxies. Technical report, Department of Computer Science, Princeton University, Jun 2012.
- [13] R. Clayton. Failures in a Hybrid Content Blocking System. In *Privacy Enhancing Technologies*, pages 78–92, Cambridge, England, 2006. Springer.
- [14] R. Clayton, S. J. Murdoch, and R. N. M. Watson. Ignoring the Great Firewall of China. In G. Danezis and P. Golle, editors, *Proceedings of the Sixth Workshop on Privacy Enhancing Technologies (PET 2006)*, pages 20–35. Springer, June 2006.
- [15] D. Bamman and B. OConnor and N. Smith. Censorship and deletion practices in Chinese social media. Online, November 2014. <http://journals.uic.edu/ojs/index.php/fm/article/view/3943/3169>.
- [16] A. Dainotti, A. Pescap, and K. Claffy. Issues and future directions in traffic classification. *IEEE Network*, 26(1):35–40, Jan 2012.
- [17] P. Dörfinger, G. Panholzer, and W. John. Entropy estimation for real-time encrypted traffic identification (short paper). In *Traffic Monitoring and Analysis*, volume 6613 of *Lecture Notes in Computer Science*, pages 164–171. Springer, 2011.
- [18] K. P. Dyer, S. E. Coull, T. Ristenpart, and T. Shrimpton. Protocol misidentification made easy with format-transforming encryption. In *Proceedings of the 20th ACM conference on Computer and Communications Security (CCS 2013)*, November 2013.
- [19] Dynamic Internet Technology Inc. DynaWeb. Proxy service. <http://dit-inc.us/dynaweb.html>. Online. October, 2014.
- [20] N. Feamster, M. Balazinska, G. Harfst, H. Balakrishnan, and D. R. Karger. Infranet: Circumventing web censorship and surveillance. In D. Boneh, editor, *USENIX Security Symposium*, pages 247–262. USENIX, 2002.
- [21] D. Fifield, G. Nakibly, and D. Boneh. Oss: Using online scanning services for censorship circumvention. In *Privacy Enhancing Technologies*, volume 7981 of *LNCS*, pages 185–204, 2013.
- [22] foe-project. <https://code.google.com/p/foe-project/>. Online. October, 2014.

- [23] J. Geddes, M. Schuchard, and N. Hopper. Cover Your ACKs: Pitfalls of Covert Channel Censorship Circumvention. In *Computer and Communications Security*, Berlin, Germany, 2013. ACM.
- [24] I. Goldberg, D. Stebila, and B. Ustaoglu. Anonymity and one-way authentication in key exchange protocols. *Designs, Codes and Cryptography*, 67(2):245–269, 2013.
- [25] A. Hintz. Fingerprinting websites using traffic analysis. In R. Dingledine and P. Syverson, editors, *Proceedings of Privacy Enhancing Technologies workshop (PET 2002)*. Springer-Verlag, LNCS 2482, April 2002.
- [26] A. Houmansadr, C. Brubaker, and V. Shmatikov. The Parrot is Dead: Observing Unobservable Network Communications. In *Symposium on Security & Privacy*, San Francisco, CA, USA, 2013. IEEE.
- [27] A. Houmansadr, T. Riedl, N. Borisov, and A. Singer. I Want my Voice to be Heard: IP over Voice-over-IP for Unobservable Censorship Circumvention. In *Proceedings of the Network and Distributed System Security Symposium - NDSS'13*. Internet Society, February 2013.
- [28] A. Houmansadr, E. L. Wong, and V. Shmatikov. No Direction Home: The True Cost of Routing Around Decoys. In *The 21<sup>st</sup> Annual Network and Distributed System Security Symposium (NDSS)*, 2014.
- [29] O. Initiative. China's Green Dam: The Implications of Government Control Encroaching on the Home PC. <https://opennet.net/chinas-green-dam-the-implications-government-control-encroaching-home-pc>.
- [30] L. Invernizzi, C. Kruegel, and G. Vigna. Message In A Bottle: Sailing Past Censorship. In *Annual Computer Security Applications Conference*, New Orleans, LA, USA, 2013. ACM.
- [31] Journalism and Media Studies Centre. Weiboscope. Online, November 2014. <http://weiboscope.jmsc.hku.hk>.
- [32] J. Karlin, D. Ellard, A. W. Jackson, C. E. Jones, G. Lauer, D. P. Mankins, and W. T. Strayer. Decoy routing: Toward unblockable internet communication. In *Proceedings of the USENIX Workshop on Free and Open Communications on the Internet (FOCI 2011)*, August 2011.
- [33] J. Karlin, D. Ellard, A. W. Jackson, C. E. Jones, G. Lauer, D. P. Mankins, and W. T. Strayer. Decoy Routing: Toward Unblockable Internet Communication. In *Free and Open Communications on the Internet*, San Francisco, CA, USA, 2011. USENIX.
- [34] S. Khattak, M. Javed, P. D. Anderson, and V. Paxson. Towards Illuminating a Censorship Monitor's Model to Facilitate Evasion. In *Free and Open Communications on the Internet*, Washington, DC, USA, 2013. USENIX.
- [35] J. Knockel, J. R. Crandall, and J. Saia. Three Researchers, Five Conjectures: An Empirical Analysis of TOM-Skype Censorship and Surveillance. In *Free and Open Communications on the Internet*, San Francisco, CA, USA, 2011. USENIX.
- [36] I7-filter. <http://research.dyn.com/2013/08/myanmar-internet/>. Online, November, 2014.
- [37] Lantern. <https://getlantern.org>. Online, November, 2014.
- [38] B. Leidl. obfuscated-openssl. <https://github.com/brl/obfuscated-openssl/blob/master/README.obfuscation>, 2010.
- [39] S. Li, M. Schliep, and N. Hopper. Facet: Streaming over videoconferencing for censorship circumvention. In *Proceedings of the 13th Workshop on Privacy in the Electronic Society, WPES '14*, pages 163–172, New York, NY, USA, 2014. ACM.
- [40] P. Lincoln, I. Mason, P. Porras, V. Yegneswaran, Z. Weinberg, J. Massar, W. A. Simpson, P. Vixie, and D. Boneh. Bootstrapping communications into an anti-censorship system. In *USENIX Workshop on Free and Open Communications on the Internet (FOCI)*, August 2012.
- [41] MailMyWeb. <http://www.mailmyweb.com>. Online, October, 2014.
- [42] W. Mazurczyk, P. Szaga, and K. Szczypiorski. Using Transcoding for Hidden Communication in IP Telephony. 2011.
- [43] meek. <https://trac.torproject.org/projects/tor/wiki/doc/meek>. Online, September, 2014.
- [44] C. MiscBrubaker, A. Houmansadr, and V. Shmatikov. CloudTransport: Using Cloud Storage for Censorship-Resistant Networking. In *Privacy Enhancing Technologies Symposium*, Springer, 2014.
- [45] H. Mohajeri Moghaddam, B. Li, M. Derakhshani, and I. Goldberg. Skypemorph: Protocol obfuscation for tor bridges. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security, CCS '12*, pages 97–108, New York, NY, USA, 2012. ACM.
- [46] nDPI. <http://www.ntop.org/products/ndpi/>. Online, November, 2014.
- [47] Nick Mathewson. <https://gitweb.torproject.org/torspec.git/blob/HEAD:/proposals/216-ntor-handshake.txt>. Online, September, 2014.
- [48] A. Panchenko, L. Niessen, A. Zinnen, and T. Engel. Website fingerprinting in onion routing based anonymization networks. In *Proceedings of the Workshop on Privacy in the Electronic Society (WPES 2011)*. ACM, October 2011.
- [49] pdfmyurl. <http://pdfmyurl.com>. Online, September, 2014.
- [50] Psiphon Inc. Psiphon. Online, November 2014. <https://psiphon.ca/en/index.html>.
- [51] Renesys. <http://research.dyn.com/2011/01/egypt-leaves-the-internet/>. Online, November, 2014.
- [52] Renesys. <http://research.dyn.com/2011/08/the-battle-for-tripolis-intern/>. Online, November, 2014.
- [53] Renesys. <http://research.dyn.com/2013/08/myanmar-internet/>. Online, November, 2014.
- [54] Renesys. <http://research.dyn.com/2013/09/internet-blackout-sudan/>. Online, November, 2014.
- [55] T. Ruffing, J. Schneider, and A. Kate. Identity-based steganography and its applications to censorship resistance, 2013. 6th Workshop on Hot Topics in Privacy Enhancing Technologies (HotPETs 2013).
- [56] J. Salowey, H. Zhou, P. Eronen, and H. Tschofenig. MessageStreamEncryption. <http://tinyurl.com/m9m17ct>, 2006.
- [57] M. Schuchard, J. Geddes, C. Thompson, and N. Hopper. Routing around decoys. In *Proceedings of the 2012 ACM Conference on Computer and Communications Security, CCS '12*, pages 85–96, New York, NY, USA, 2012. ACM.
- [58] R. Smits, D. Jain, S. Pidcock, I. Goldberg, and U. Hengartner. Spator: Improving tor bridges with single packet authorization.
- [59] R. Smits, D. Jain, S. Pidcock, I. Goldberg, and U. Hengartner. Bridgespa: Improving tor bridges with single packet authorization. In *Proceedings of the 10th Annual ACM Workshop on Privacy in the Electronic Society, WPES '11*, pages 93–102, New York, NY, USA, 2011. ACM.
- [60] Snort. <https://www.snort.org>. Online, November, 2014.
- [61] R. Sommer and V. Paxson. Outside the closed world: On using machine learning for network intrusion detection. In *In Proceedings of the IEEE Symposium on Security and Privacy*, 2010.
- [62] Steven J. Murdoch. Pluggable Transport Component Architecture. <https://gitweb.torproject.org/sjm217/torspec.git/tree/pt-components.txt?h=pt-components>.
- [63] Q. Sun, D. R. Simon, Y.-M. Wang, W. Russell, V. N. Padmanabhan, and L. Qiu. Statistical identification of encrypted web browsing traffic. In *IEEE Symposium on Security and Privacy*, May 2002.
- [64] T. Wilde. Great firewall tor probing. Online, November 2014. <https://gist.github.com/da3c7a9af01d74cd7de7>.
- [65] The Tor Project. Fog. Online, November 2014. <https://gitweb.torproject.org/pluggable-transport/fog.git>.
- [66] Tor. <https://github.com/Yawning/obfs4/blob/5bdc376e2abaf5ac87816b763f5b26e314ee9536/doc/obfs4-spec.txt>. Online, September, 2014.
- [67] Tor. <https://gitweb.torproject.org/pluggable-transport/obfsproxy.git/blob/HEAD:/doc/obfs2/obfs2-protocol-spec.txt>. Online, September, 2014.
- [68] Tor. <https://gitweb.torproject.org/pluggable-transport/obfsproxy.git/blob/HEAD:/doc/obfs3/obfs3-protocol-spec.txt>. Online, September, 2014.
- [69] Tor. <https://www.torproject.org/docs/pluggable-transport.html.en>. Online, September, 2014.
- [70] Tor. Online, November 2014. <https://www.torproject.org>.
- [71] M. C. Tschantz, S. Afroz, V. Paxson, and J. D. Tygar. On Modeling the Costs of Censorship. 2014.
- [72] Ultrasurf. <http://ultrasurf.us>. Online, October, 2014.
- [73] E. Y. Vasserman, N. Hopper, and J. Tyra. Silent knock : practical, provably undetectable authentication. *Int. J. Inf. Sec.*, 8(2):121–135, 2009.
- [74] B. Wiley. Blocking-resistant protocol classification using bayesian model selection. Technical report, University of Texas at Austin, 2011.
- [75] B. Wiley. Dust: A blocking-resistant internet transport protocol. Technical report, School of Information, University of Texas at Austin, 2011.
- [76] P. Winter and S. Lindskog. How the Great Firewall of China is Blocking Tor. In *Free and Open Communications on the Internet*, Bellevue, WA, USA, 2012. USENIX.
- [77] P. Winter, T. Pulls, and J. Fuss. ScrambleSuit: A Polymorphic Network Protocol to Circumvent Censorship. In *Workshop on Privacy in the Electronic Society*, Berlin, Germany, 2013. ACM.
- [78] C. V. Wright, F. Monrose, and G. M. Masson. On inferring application protocol behaviors in encrypted network traffic. *J. Mach. Learn. Res.*, 7:2745–2769, Dec. 2006.
- [79] E. Wustrow, C. M. Swanson, and J. A. Halderman. Tapdance: End-to-middle anticensorship without flow blocking. In *23rd USENIX Security Symposium (USENIX Security 14)*, pages 159–174, San Diego, CA, Aug. 2014. USENIX Association.
- [80] E. Wustrow, S. Wolchok, I. Goldberg, and J. A. Halderman. Telex: Anticensorship in the Network Infrastructure. In *USENIX Security Symposium*, San Francisco, CA, USA, 2011. USENIX.

- [81] W. Zhou, A. Houmansadr, M. Caesar, and N. Borisov. SWEET: Serving the Web by Exploiting Email Tunnels. In *Hot Topics in Privacy Enhancing Technologies*, Bloomington, IN, USA, 2013. Springer.
- [82] T. Zhu, D. Phipps, A. Pridgen, J. R. Crandall, and D. S. Wallach. The velocity of censorship: High-fidelity detection of microblog post deletions. In *USENIX Security*, pages 227–240. USENIX, 2013.