

De-anonymizing Encrypted Video Streams

Master's Thesis

7 September 2018

Stefano Peverelli

pstefano@student.ethz.ch

supervised by

Melissa Licciardello

Systems Group
Department of Computer Science
ETH Zürich

Abstract

In the last recent years streaming services such as Netflix, Youtube, Amazon Prime Video, Hulu and others, have become the main source for video content delivery to the public. With the effort of private companies and of the AOM consortium [1], various coding formats and streaming techniques have been refined and have been vastly adopted in the industry. *Variable Bitrate Encoding* and *Adaptive Bitrate Streaming*, between others, are two techniques that increasingly often get coupled together to improve high quality streaming of media content over HTTP.

Netflix relies on first encoding a title at a *variable* bitrate, and streaming it using DASH *Dynamic Adaptive Streaming over HTTP*, an instance of Adaptive Bitrate Streaming originally developed by MPEG. In DASH each media file gets encoded at multiple bitrates, which are then partitioned into smaller segments and delivered to the user over HTTP. As of today, Netflix makes extensive use of DASH, by encoding each title in their catalogue at different quality levels, so that the user is served tailored-size content depending on network conditions, and its device capabilities [2].

DASH and *variable-bitrate-encoding*, have proven to be leaking potentially harmful data. Reed et Al. [3] have shown, how despite a recent upgrade in Netflix infrastructure to provide HTTPS encryption to video traffic, it is possible to recover unique fingerprints for each title, and to use them to match against captured traffic. In their study, they make use of adudump [4] a command-line program that can run on passive TAP devices [5] or on a live network interface, and uses TCP and ACKs sequences to infer the size of application data units *ADUs* transferred over each TCP connection in real time.

Our approach reiterates parts of Reed et Al.'s work, but cannot rely on their every assumption and discovery, due to constant changes that Netflix is integrating into their encoding and streaming algorithms. Moreover our intent is focused on finding out if, by analyzing coarse-grained traffic data, we are able to identify a video based on its *bitrate-ladder*.

TODO: refine it and add results

Contents

1	Introduction	7
1.1	Motivation	7
1.1.1	Per-Title Encoding	7
1.1.2	User's Privacy	9
1.2	Related Work	10
1.3	Main Objective	11
1.4	Structure of this Report	11
2	ISPs as adversaries	13
2.1	Attack Scenarios	13
2.2	Video Fingerprinting	14
2.2.1	The a-b-t model	14
2.2.2	Inference	15
2.2.3	Identifying Netflix traffic	15
2.3	Capturing video traffic	17
3	Approach	19
	Bibliography	21

Introduction

According to the latest Cisco's VNI [6], video will account for 82% of all IP traffic in Europe by 2021; in addition, the overall IP traffic per person will triplicate from 13GB to 35GB. These forecasts clearly picture the growth of the streaming industry, posing, at the same time an important question on the present and future states of the final user's privacy.

As shown by Reed et Al. [3] anonymity of user's viewing activity is at risk. Not for the use that Netflix or other streaming services do of user's session data, but because of the risk of a man-in-the-middle attack *MITM* carried by an *evil* party that has control over the flow of packets over a network.

In particular, they have shown how the adoption of HTTPS to protect video streams from Netflix *CDNs* to user's end devices, does not hold against passive traffic analysis.

1.1 Motivation

The goal of this project is to replicate part of the work conducted by Reed et Al. and to investigate the possibility of identifying a Netflix stream solely based on the observed average bandwidth. This, follows from the intuition that *per-title encoding* embeds the nature and the complexity of video frames in a unique way, that may reveal the identity of the content being streamed.

1.1.1 Per-Title Encoding

In December 2015 Netflix announced [2] that it was introducing a new method to analyze the complexity of each title and find the best encoding recipe based on it. Their goal with the adoption of per-title encoding was to provide users with better quality streams at a lower bandwidth.

Chapter 1 Introduction

Before then, each title was encoded with a *Fixed Bitrate Ladder*; their pipeline returned a list of $\{\text{Bitrate}, \text{Resolution}\}$ pairs that represented the sufficient bitrate to encode the stream at a certain resolution (Table 1.1), with no visible artifacts.

Bitrate (kbps)	Resolution
235	320×240
375	384×288
560	512×384
750	512×384
1050	640×480
1750	720×480
2350	1280×720
3000	1280×720
4300	1920×1080
5800	1920×1080

Table 1.1: Netflix original's Fixed Bitrate Ladder.

This "one-size-fits-all" ladder, as reported, achieved good results in the encoded video's perceived quality (PSNR [7]) given the bitrate constraint, but, would not perform optimally under certain conditions. For instance, high detailed scenes with sudden changes of light, or rapid transitions of camera shots, would require more than 5800kbps ; in contrast, more static frames, as in animated cartoons, may be encoded at higher resolutions maintaining the same bitrate level.

In summary they noticed how in certain cases, the produced encoding would either present some small artifacts (*e.g.* complex scenes), or waste bandwidth, (*e.g.* static, plain scenes). For this reason, they came up with per-title encoding.

In order to find the best fitting bitrate ladder for a particular title, there are several criterias that they took into account, the principal ones being:

- How many quality levels should be encoded to obtain a *JND* between each of them.
- Best $\{\text{Resolution}, \text{Bitrate}\}$ pair for each quality level
- Highest bitrate required to achieve the best perceivable quality

As aforementioned, each title's perceived video quality, gets computed as a measure of *Peak signal-to-noise ratio*. The comparison is performed between the produced encode, upsampled to $1080p$, and the original title in $1080p$, and the best $\{\text{Bitrate}, \text{Resolution}\}$ pair is assigned to that specific quality level, as depicted in Table 1.2.

In Figure 1.1, we can see the impact of per-title encoding on the original bitrate ladder: in order to achieve the same perceivable quality level (point **B** and **C**), it requires a lower bitrate to be encoded to (point **A**). Moreover, with around the same bitrate, one can see

Resolutions	Fixed Bitrate Ladder (kpbs)	Per-Title Bitrate Ladder (kpbs)
320×240	235	150
384×288	375	200
512×384	560	290
512×384	750	
640×480	1050	
720×480	1750	440
720×480		590
1280×720	2350	830
1920×1080	3000	1150
1920×1080	4300	1470
1920×1080	5800	2150
1920×1080		3840

Table 1.2: Comparison between the two different approaches for the same title: note how different titles may have different numbers of quality levels. For each movie, the minimum number of quality levels gets computed to produce a just-noticeable-difference (JND), when switching bitrates during playback.

how per-title encoding can achieve a higher resolution compared the fixed case (point **A** and **D** respectively). It follows obviously that, holding to a high-quality stream while maintaining or lowering the used bandwidth is key: the end user will get same or better quality then before, at a lower bandwidth.

1.1.2 User's Privacy

As of 2018, data began the most valuable commodity on the planet [8], and with its value rising, society starts to question how to legislate to protect both industry's and user's rights.

In 2016, Netflix announced the introduction of HTTPS to protect the content being streamed to users. With the addition of TLS on top of HTTP, Netflix aim, was to avoid the risk of eavesdropping on unsecure connections, protecting themselves and users from third-party applications and governments potentially collecting viewer's data and streaming habits.

Avoiding deep packet inspection from potential eavesdropper certainly adds another layer of security, but given the narrow relationship between Netflix and ISPs, one might conclude that the chance that IP packets do not get inspected by ISPs (especially in the United States, after the reclassification of ISPs as Title I *Information Services* [9]), is still a matter of mere trust given to the platform-provider relationship.

Cases of user's data breaches as the Kanopy one [10], are a testimony of how easy would be for an attacker to extract information to the point at which users could become

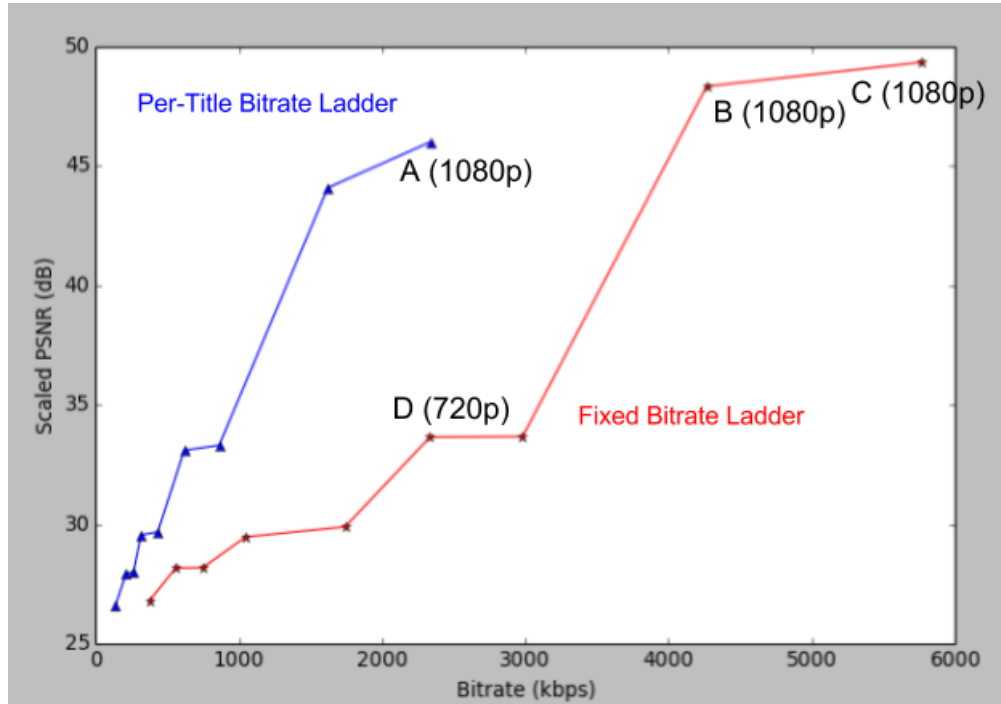


Figure 1.1: Difference between per-title vs. fixed bitrate ladders.

identifiable by the solely information the platform was collecting in their internal log files. The content of which, include between others: timestamps, geo-location data, client-device informations and IP addresses.

The ability to cluster people based on just their video-streaming habits, poses a potential threat for how the data could be processed and used by parties with access to it. One could imagine how government agencies could easily get in possess of sensitive information about the nature of the content a particular user is interested about, or how ISPs could profit from selling data profiles to advertisement companies that in turn would exploit their information to improve per-user recommendation algorithms.

Considering this trend, it is for us crucial to investigate how parties that can have access to transport-layer information, (more details in Chapter 3), could exploit per-title encoding to identify video traffic.

1.2 Related Work

As previously mentioned, our work is mainly inspired by Reed et Al.'s [3] research paper, in which they presented a novel method to de-anonymize encrypted netflix stream in real-time with limited hardware requirements. Their system was able to identify a video using uniquely TCP/IP headers, by making use of *adudump*, a program built on

top of *libpcap* [11], a powerful C library for network traffic capture. They acquired for each video, metadata information with a tampering tool, and then matched adudump traces against with. The evaluation of their method revealed that they could identify majority of videos by recording only 20 minutes of traffic each.

An earlier proof on how bandwidth analysis could reveal the content of encrypted traffic, was given by Saponas et Al. [12], in their "SlingBox-Pro" case study, that exposed how, by recording network traffic and producing and combining different trace levels, they were able to identify 98% of 40 minutes video traces.

Similar work has been conducted also by Moser [13], who analyzed how bitrate ladders could uniquely embed the identity of Netflix titles. In his work he further studies the impact of the aggregation of each video's segment bandwidth on the overall accuracy of his system.

1.3 Main Objective

The main goal of this project is to build a system that can manipulate network bandwidth and observe video traffic from Netflix, to verify the intuition that each title could be identified by its own bitrate ladder.

Furthermore we investigate and discuss possible countermeasures that streaming providers could adopt to preserve users privacy. A detailed explanation of our approach is presented in Chapter 3

1.4 Structure of this Report

In this section, we outline of the contents and the structure of this report.

Chapter 1 introduced our motivation behind the project, presented featured work that influenced and inspire our approach, and it outlined the goals we would like to ultimately achieve.

Chapter 2 presents the attack scenario from an ISP perspective, (the ISP acting as the adversary), the infrastructure needed to acquire the data, the nature of the information that could be inferred, and the consequences that might arise.

Chapter 3 shows our version of the attack, in a different context, but up to some extent, with similar conditions to the one depicted in Chapter 2. It also highlights the similarities and differences from the featured approaches we followed.

In Chapter 4 we evaluate our system, we present results, and discuss the relevance of such a method.

Chapter 5 tries to summarize and draw conclusion based on the claims made at the beginning of the project.

ISPs as adversaries

Of great importance to our approach, is acquiring real-time network traffic, with downstream throughput being our only focus; here discuss the potential scenarios that allow adversaries to obtain such information, and describe how, in particular ISPs, could make use of it, to identify video streaming and profile users based on their streaming habits.

2.1 Attack Scenarios

Given the nature of modern Internet infrastructure, an adversary interested in eavesdropping a particular communication, only needs to compromise a node on the path the communication travels through. An *on-path* attacker could easily gain passive access to network and transport layers, and start capturing network traffic. This can include malicious or compromised Wi-Fi access points, routers, tapped network cables and ISPs.

Rather than just attacking physical devices, leaks of information could occur in network connections, in which an attacker physically close to the victim, could make use of a *Wi-Fi sniffer* to estimate traffic by capturing physical layer WLAN packets. Information may be encrypted by 802.11, and the sniffer may not take into account packet retransmissions at the session layer nor multiple TCP/IP flows on the same link, potentially causing noise on the observation. Despite so, Reed et Al. [14], have shown that it is still possible to estimate WAP-to-client throughput and use it to identify the content being streamed.

In addition to the above, *side-channel eavesdropping* can exploit information about the network structure, to saturate a link between the user and the server, and estimate fluctuations of congestion by sending probes remotely and observing queueing delays in routers [15].

We will now present the relevant phases and tools needed for an ISP to perform such attack.

2.2 Video Fingerprinting

Assuming that the ISP has no direct access to the CDN titles are stored in, it needs to build a database of video traffic to match against future video stream captures. To build such a database, the ISP needs to have access to a network interface, control its inbound bandwidth, (to get different levels of quality for each title), and capture incoming traffic passing through it.

In order to control the incoming bandwidth, the ISP could either limit it directly onto a generic L4 switch (probably more stable), or decide to connect any UNIX-like machine and throttle the throughput of its main ethernet interface. We will consider the scenario in which the ISP limits the bandwidth of the ethernet interface of the switch.

The value of the enforced bandwidth determines the quality (bitrate) of the content that will be captured. Assuming that each title has a unique bitrate ladder, the ISP should come up with an ad-hoc policy for every video to faithfully reconstruct the quality levels of it. While correct, a more viable approach to this problem, would be to just consider a range of bitrates capable of spanning the space of possible quality levels.

In our own version of the attack, presented in Chapter 3, we use the values in Table 2.1.

Bandwidth levels (Mbps)												
0.6	0.8	1.2	2	3.5	4.2	4.8	5.5	6.5	7.05	10	15	20

Table 2.1: Enforced bandwidth levels.

The ISP can now connect a UNIX-like machine to the switch’s network interface with bandwidth limit b , and invoke `adudump` to infer the size of each *application data unit* ADU by processing TCP/IP packet header traces that generate *a-b-t* connection vectors [16].

2.2.1 The *a-b-t* model

The lifetime of a TCP endpoint that sends and receive data units is not only dictated by the time spent on these operations, but also by quiet times in which the TCP connection remains idle, waiting for upper layers to formulate new demands. Clearly, longer lifetimes may have a huge impact overall, due to the fact that resources needed to handle TCP state, remain reserved for a longer time. In contrast, ADUs that are sent within a short period of time, get aggregated by the TCP window mechanism, reducing the exchange of requests/responses between source and destination.

These ideas have been formalized into the *a-b-t model*, that describes TCP connections as sets of ADU exchanges and quiet times. The term *a-b-t* is descriptive of the building blocks of the model: *a-type* ADUs, sent from the connection initiator to the connection

acceptor, *b-type* ADUs, sent from the responder back to the initiator, and quiet times t' 's, during no data segments are exchanged.

The model has two different flavors depending on whether ADU interleaving is sequential or concurrent: the former is used for modeling connections for which only one ADU is sent from one endpoint to the other at a given point in time, (one endpoint will wait until the other has completed the transmission of the current ADU before sending a new one), while the latter is used for modeling connections in which both endpoints send and receive ADUs simultaneously.

As one can notice from a sample traffic trace from Netflix, the two endpoints behave in a sequential fashion, as shown by the SEQ flag in Listing 2.1.

```
ADU: 1565880580.152388 192.168.0.157.54924 <1 45.57.19.134.443 198837 SEQ 0.804088
ADU: 1565880581.321246 192.168.0.157.54924 <1 45.57.19.134.443 198760 SEQ 0.907568
ADU: 1565880582.688335 192.168.0.157.54924 <1 45.57.19.134.443 199407 SEQ 1.052862
ADU: 1565880591.912352 192.168.0.157.54974 <1 45.57.19.134.443 525046 SEQ 0.058709
ADU: 1565880592.168340 192.168.0.157.54968 <1 45.57.19.134.443 275568 SEQ 0.558654
ADU: 1565880592.168344 192.168.0.157.54974 <1 45.57.19.134.443 305115 SEQ 0.062028
ADU: 1565880592.172495 192.168.0.157.54982 <1 45.57.19.134.443 198735 SEQ 0.670224
ADU: 1565880592.697314 192.168.0.157.54974 <1 45.57.19.134.443 236474 SEQ 0.179115
ADU: 1565880592.820292 192.168.0.157.54982 <1 45.57.19.134.443 286547 SEQ 0.106913
ADU: 1565880592.820395 192.168.0.157.54968 <1 45.57.19.134.443 198530 SEQ 0.310506
```

Listing 2.1: Incoming traffic trace of Mulan, captured at 10Mbps (first 10 segments).

2.2.2 Inference

Figure 2.1 shows an example of how Adudump's inference process works. The connection starts when the three-way handshake completes, (marked as a SEQ record). The monitor then records the time the first data segment is sent from the client. The server then acknowledges the previous request, and adudump infers the size of the completed ADU to be 372 bytes, generating a record with the ADU's size direction and think-time. Note how adudump generates no record until it infers that the ADU is complete, moreover the think times reported, are relative to the position of the monitor in the network.

2.2.3 Identifying Netflix traffic

In order to identify specific video traffic from a Netflix CDN, the ISP could either inspect the DNS response packets or TLS handshake messages, and perform a DNS lookup for each video every time, or to maintain an updated list of IP addresses with mappings to the recorded traffic. Then all the ISP is left with, is a mixed trace of audio and video segments. Each Netflix audio segments represents about 16 seconds of playback, while every video segment corresponds to a 4 second playback. The ISP could now decide to either filter out audio segments by either thresholding on the size of ADUs, as done by

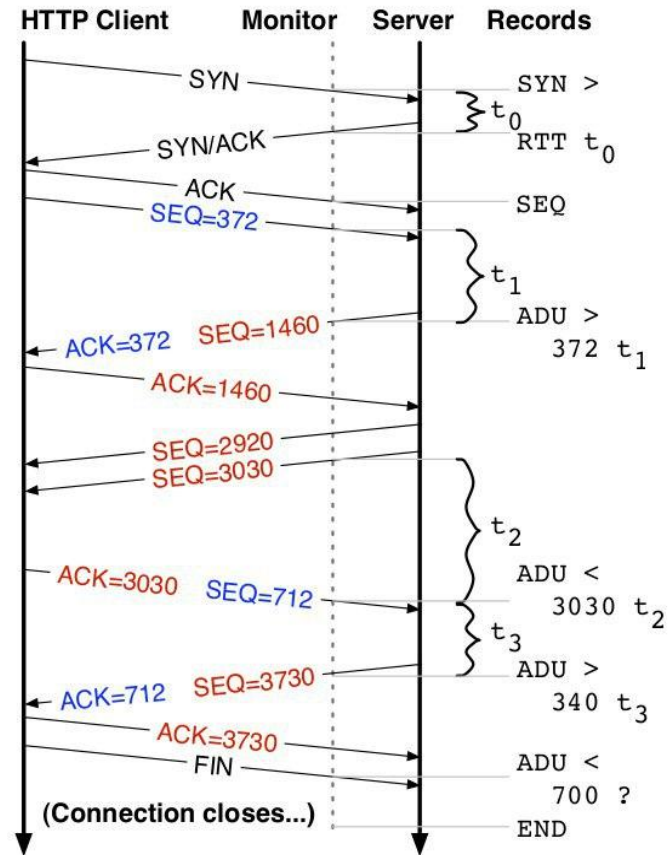


Figure 2.1: Detailed segment size inference in adudump.

Reed et Al (at 200 kilobytes), or to consider audio as part of the fingerprint of a title. The 200 kilobytes threshold represent a good heuristic when considering traces recorded at high bandwidth levels, in fact, audio segments sizes are around 198 kilobytes, whereas if considering the same title at a low level of bandwidth (less than 1 Mbps indicatively), one cannot distinguish audio and video segments by their size. For this reason, we consider an approach that does not filter out audio segments.

2.3 Capturing video traffic

This, requires the ISP to be in possess of a generic L4 switch that can mirror traffic from a port to another one. Then any UNIX-like system that implements *libpcap* is suitable for capturing inbound traffic on the mirrored port, as presented below in Figure 2.2.

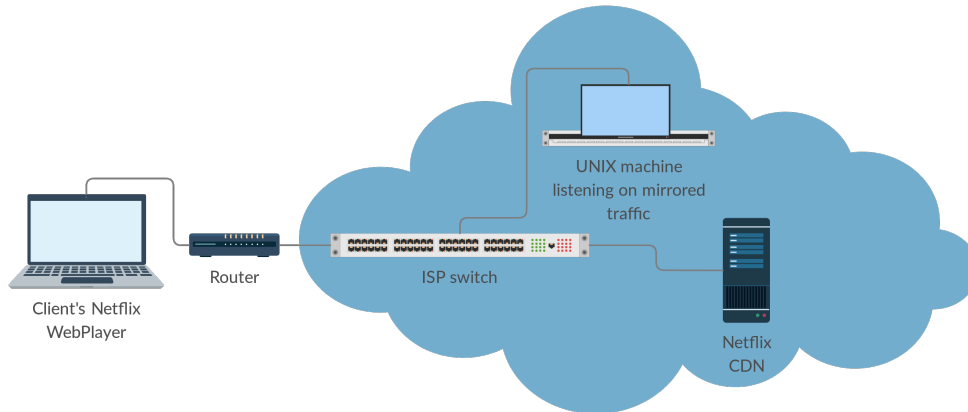


Figure 2.2: Traffic capture scenario.

Approach

After having shown how an ISP could exploit *per-title-encoding* to reconstruct video traffic traces, we will now present our contribution and practical approach for building a database of fingerprints to match and possibly identify future captured traffic against.

Bibliography

- [1] Wikipedia. Alliance for open media, 2019. URL https://en.wikipedia.org/wiki/Alliance_for_Open_Media.
- [2] Netflix TechBlog. Per-title encode optimization, 2015. URL <https://medium.com/netflix-techblog/per-title-encode-optimization-7e99442b62a2>.
- [3] Andrew Reed and Michael Kranch. Identifying https-protected netflix videos in real-time. In *Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy, CODASPY '17*, pages 361–368, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4523-1. doi: 10.1145/3029806.3029821. URL <http://doi.acm.org/10.1145/3029806.3029821>.
- [4] Jeff Terrell, Kevin Jeffay, F. Donelson Smith, Jim Gogan, and Joni Keller. Passive, streaming inference of tcp connection structure for network server management. IEEE International Traffic Monitoring and Analysis Workshop, 2009.
- [5] Wikipedia. Network tap, 2019. URL https://en.wikipedia.org/wiki/Network_tap.
- [6] Business Insider Intelligence. Video will account for an overwhelming majority of internet traffic by 2021, 2019. URL <https://www.businessinsider.com/heres-how-much-ip-traffic-will-be-video-by-2021-2017-6?r=US&IR=T>.
- [7] Wikipedia. Peak signal-to-noise ratio, 2019. URL https://en.wikipedia.org/wiki/Peak_signal-to-noise_ratio.
- [8] The Economist. The world’s most valuable resource is no longer oil, but data, 2018. URL <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>.
- [9] Wikipedia. Net neutrality in the united states, 2019. URL https://en.wikipedia.org/wiki/Net_neutrality_in_the_United_States.
- [10] @xxdesmus. Kanopy.com leaking api and website access logs, 2019. URL <https://rainbowtabl.es/2019/03/21/kanopy-data-leak/>.
- [11] tcpdump. URL <https://www.tcpdump.org/>.

Bibliography

- [12] Scott Saponas, Jonathan Lester, Carl Hartung, Sameer Agarwal, and Tadayoshi Kohno. Devices that tell on you: Privacy trends in consumer ubiquitous computing. In *Proceedings of the 16th USENIX Security Symposium*. USENIX, August 2007. URL <https://www.microsoft.com/en-us/research/publication/devices-tell-privacy-trends-consumer-ubiquitous-computing/>.
- [13] Florian Moser. Identifying encrypted online video streams using bitrate profiles, 2018. URL <https://github.com/famoser/network-experiments>.
- [14] A. Reed and B. Klimkowski. Leaky streams: Identifying variable bitrate dash videos streamed over encrypted 802.11n connections. In *2016 13th IEEE Annual Consumer Communications Networking Conference (CCNC)*, pages 1107–1112, Jan 2016. doi: 10.1109/CCNC.2016.7444944.
- [15] Xun Gong, Nikita Borisov, Negar Kiyavash, and Nabil Schear. Website detection using remote traffic analysis. In Simone Fischer-Hübner and Matthew Wright, editors, *Privacy Enhancing Technologies*, pages 58–78, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. ISBN 978-3-642-31680-7.
- [16] F Donelson Smith, Felix Hernández-Campos, Kevin Jeffay, and David Ott. What tcp/ip protocol headers can tell us about the web. volume 29, pages 245–256, 06 2001. doi: 10.1145/378420.378789.



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

De-anonymizing encrypted video streams

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

Stefano

First name(s):

Peeverelli

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Zurich, 07.09.2019

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.