

# Do Security Reports Meet Usability?

Lessons Learned from Using Actionable Mitigations for Patching TLS Misconfigurations

Salvatore Manfredi

Security & Trust, FBK  
Trento, Italy  
DIBRIS, University of Genoa  
Genoa, Italy  
smanfredi@fbk.eu

Giada Sciarretta

Security & Trust, FBK  
Trento, Italy  
giada.sciarretta@fbk.eu

Mariano Ceccato

Department of Computer Science, University of Verona  
Verona, Italy  
mariano.ceccato@univr.it

Silvio Ranise

Security & Trust, FBK  
Trento, Italy  
Department of Mathematics, University of Trento  
Trento, Italy  
ranise@fbk.eu

## ABSTRACT

Several automated tools have been proposed to detect vulnerabilities. These tools are mainly evaluated in terms of their accuracy in detecting vulnerabilities, but the evaluation of their usability is a commonly neglected topic. Usability of automated security tools is particularly crucial when dealing with problems of cryptographic protocols for which even small—apparently insignificant—changes in their configuration can result in vulnerabilities that, if exploited, pave the way to attacks with dramatic consequences for the confidentiality and integrity of exchanged messages. This becomes even more acute when considering such ubiquitous protocols as the one for Transport Layer Security (TLS for short). In this paper, we present the design and the lessons learned of a user study, meant to compare two different approaches when reporting misconfigurations. Results reveal that including contextualized actionable mitigations in security reports significantly impact the accuracy and the time needed to patch TLS vulnerabilities. Along with the lessons learned, we share the experimental material that can be used during cybersecurity labs to let students configure and patch TLS first-hand.

## CCS CONCEPTS

• **Networks** → **Network security**; • **Security and privacy** → *Social aspects of security and privacy.*

## KEYWORDS

vulnerability detection, usability study, actionable mitigations, security reports, TLS misconfiguration

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*ARES 2021, August 17–20, 2021, Vienna, Austria*

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-9051-4/21/08...\$15.00

<https://doi.org/10.1145/3465481.3469187>

## ACM Reference Format:

Salvatore Manfredi, Mariano Ceccato, Giada Sciarretta, and Silvio Ranise. 2021. Do Security Reports Meet Usability?: Lessons Learned from Using Actionable Mitigations for Patching TLS Misconfigurations. In *The 16th International Conference on Availability, Reliability and Security (ARES 2021), August 17–20, 2021, Vienna, Austria*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3465481.3469187>

## 1 INTRODUCTION

Transport Layer Security (TLS) consists of a set of cryptographic protocols designed to provide secure communications over a network. TLS is widely used in client-server applications to secure all the communications by preventing eavesdropping and tampering. It is mainly used to secure the traffic between a website and a web browser (HTTPS protocol). Another use of TLS is on top of Transport Layer protocols such as File Transfer Protocol (FTP) for the transfer of computer files or Simple Mail Transfer Protocol (SMTP) for electronic mail transmission.

The TLS deployment process requires a non-trivial amount of knowledge [41, 46] as this results in having to correctly set a server certificate, choose the available protocols, ciphers and enabling other security mechanisms (e.g., HSTS). Its popularity has encouraged attackers to find vulnerabilities. The types of attacks vary widely and include the renegotiation of cipher suites to exploit weak encryption algorithms [39], the knowledge of initialization vectors to retrieve symmetric keys [24], and the use of libraries to exploit poor certificate validation in deployments where clients are non-browsers [20]. The downside of allowing a high level of customization transfers on the system administrators the burden of securing the deployments.<sup>1</sup> However, according to [47], security does not get enough attention as more than 70,000 of the most popular websites [55] still support officially deprecated versions of the TLS protocol [31].

To help administrators in deploying secure TLS instances, during the years, a variety of tools have been developed to assist system administrators and ease their work. These tools are usually able to

<sup>1</sup><https://acloudguru.com/blog/engineering/cloud-governance-and-managing-risk>

verify TLS implementations, analyzing the deployment and providing a list of the possible attacks but miss two important features: (i) an *explanation* of the identified vulnerabilities and the attacks that may exploit them, and (ii) *actionable hints* on how to mitigate the attacks. The *explanation* would allow administrators to learn about the attacks and put the potential security problems in the right context (such as impact and likelihood), help them to distinguish relevant from irrelevant information when searching for more details, and even ease the acceptance of *actionable hints* that can even take the form of code snippet to be copy-pasted in a configuration file. Without these two features, a system administrator is left with the burden of finding enough information to define an appropriate mitigation strategy. This task is far from trivial as this kind of information is spread across several sources ranging from scientific papers to blog posts; each one with its jargon and background assumptions. Even disregarding the effort to collect enough material to mitigate a security problem, administrators should have enough skills to understand the (often subtle) details and turn the information in a concrete strategy to fix the problem.

We think that a tool should go beyond than just identifying vulnerabilities, and that contextualized actionable suggestions would be very effective in practically explaining system administrators how to fix a wrong configuration and remove a specific security problem. In this paper, we present an experimental assessment of this hypothesis by quantifying the benefit of providing mitigation hints on the capability of system administrators to (correctly and quickly) patch a TLS misconfiguration.

During the experiment, bachelor and master students were asked to play the role of unexperienced system administrators who should patch defective TLS configuration files.

The findings prove the effectiveness of security reports containing high level descriptions of contextual information about identified security problems (e.g., of the identified vulnerabilities) together with actionable mitigations as even un-experienced users were able to successfully mitigate complex attacks.

Our contributions can be summarized as follows:

- we formally validate the hypothesis that reports containing actionable hints are more efficient as we learned that they drastically decrease the time required to patch TLS vulnerabilities and improve both problem identification and its resolution;
- we share the vulnerable VMs, slides, questionnaires as they can be used both to replicate our study and as valuable asset for educational scenarios.

*Plan of the paper.* Section 2 provides the necessary background notions on TLS and its vulnerabilities. Section 3 contains an overview of the state-of-the-art. Section 4 describes the decision process that led to the choice of the TLS analyzer used an experimental user study involving Bachelor and Master degree students. Section 5 describes the experimental framework that we defined to collect empirical evidence on the reports' effectiveness and Section 6 reports the results of this study. These results are then discussed in Section 7 in terms of their impact and implications. Section 8 concludes the paper and highlights future work.

## 2 BACKGROUND

We provide some background notions to better understand the experimentation. We briefly describe the general structure of the TLS protocol in Section 2.1 and give a concise description of two known vulnerabilities in Section 2.2.

### 2.1 Transport Layer Security

TLS has been designed to provide both confidentiality and integrity between communicating entities [17] and is composed of two layers: the Handshake and the Record protocol.

The Handshake protocol can provide either mutual or one-way authentication and allows the parties to exchange all the information required to establish a reliable session, this includes the choice of a set of algorithms that will be used. This set, called cipher suite, is proposed by the client within the first handshake message and then chosen by the server. The set of cipher suites supported by the server is chosen by the system administrator who deployed it.

The Record protocol is deployed on top of a transport protocol (such as TCP) and encapsulates the messages coming from higher levels, it ensures confidentiality by using symmetric encryption algorithms and integrity by calculating the hash of the messages being sent. The keys and the algorithms used are the ones agreed during the handshake.

### 2.2 Vulnerabilities on TLS

On protocol versions prior to 1.3, TLS suffers from a wide set of vulnerabilities [56]. While some of the attacks are due to flaws in the logic of the protocol, others exploit the support of now-deprecated cipher suites or (in)voluntary weakening of security properties to bypass the authentication process (such as accepting self-signed certificates [45]). In this paper, we focus on the two used in our experiments (see Section 5.2): CRIME and BREACH. Both are compression side-channel attacks that combine three elements: the presence of a recurring part within the transmitted messages, the fact that both TLS and HTTP do not hide the length of the sent data and the availability of DEFLATE [27], a compression algorithm that reduces the size of an input by replacing duplicate strings with a reference to their last occurrence.

CRIME [44] is a security exploit that allows an attacker to decrypt the transmission if the parties agreed to use TLS-level compression. Supposing the attacker's intention is to steal a session cookie (whose ownership authenticates the user), the attack is performed by creating different client's messages containing guesses. Due to the DEFLATE usage, if the guess is wrong, the size of the server's response will be bigger than a valid one.

BREACH [22] exploits the same mechanism but using HTTP-level compression. Once the attacker has correctly guessed the first character of the shared secret, it starts a phase of trial-and-error in which she/he increasingly guesses bigger parts of the secret until its completion.

Being optional, the presence of DEFLATE can be seen as the single point-of-failure and thus its deactivation – that in Apache consists in changing two different files – is the suggested mitigation for both attacks.

**Table 1: TLS Analyzers - Features Comparison**

Features	ssl-enum-ciphers	ssllscan	SSL Server Test	sslyze	testssl.sh	TLSSLed
Open source & downloadable	✓	✓	X	✓	✓	✓
Actively maintained	●○○	●●●	●●●	●●●	●●●	●○○
TLS vulnerability checks	●○○	●○○	●●●	●○○	●●●	●○○
Standalone	●○○	●○○	●●●	●○○	●●●	●○○
Highly customizable scan	●○○	●●○	●○○	●○○	●●●	●○○

### 3 RELATED WORK

*Usability Studies in Cyber Security.* Several usability studies have been conducted in order to assess tools and methodologies in different cyber-security domains, such as password storage, penetration testing and code obfuscation.

Naiakshina et al. conducted qualitative usability studies either with students [43] and with freelance developers working remotely [42], asking them to build a password storage mechanism. These studies showed that participant security knowledge does not guarantee the delivery of secure software.

In the domain of risk assessment, Allodi et al. measured the accuracy [4] and the difficulty [3] for students (with different technical education) in using the Common Vulnerability Scoring System to assess the severity of software vulnerabilities. Labunets et al. conducted a series of empirical evaluations to compare the effectiveness of two classes of threats-analysis methods [37] and the comprehensibility of two risk model representations [36].

Scandariato et al. conducted a series of controlled experiments to compare static analysis and penetration testing tools, in terms of how well they support developers in accurately detecting vulnerabilities [52], and then in fixing the code [12].

Ceccato et al. measured how code obfuscation influences the correctness and effectiveness of understanding and change tasks [10]. In a successive work, the same authors presented an extension with a larger set of experiments conducted on more obfuscation techniques [8]. Then, their replication package has been used by Hänsch et al. [26] to conduct a similar experiment and assess a slightly different set of obfuscations. Viticchié et al. empirically evaluated the attack delay introduced by a data obfuscation [64] and by code splitting [63]. They confirmed that attacks are still possible on protected programs, but they are delayed by a factor of six for that technique.

Compared to the described studies, our focus is different: we tried to understand how a sysadmin could be guided toward a correct configuration of vulnerable webservers using actionable hints.

*Impact of providing hint suggestions.* The following articles focus instead on how awareness and documentation affect the usage and related maintenance of specific technologies; they start from hypothesis similar to ours but focus on the difficulty rather than proposing a solution.

Acar et al. [2] performed a systematic investigation on how the documentation available to developers directly affects both security and privacy properties, finding that most developers do use search engines and StackOverflow to address issues, leading to poor implementation results. Gorski et al. [23] evaluated the impact

of providing security advice in case of API misuse, proving that the offered advices impacted positively on code security and did not affect the overall usability of the interface.

Krombholz et al. investigated the mental models of both users and sysadmins, their results show a large amount of misconceptions about threat models, protocol components and also the very same benefits of using HTTPS [33]. In [34], they also performed a series of controlled experiments to highlight the difficulties of deploying HTTPS, proving that it is far too complex even for people with adequate expertise. The findings of the latter have been partially verified by Bernhard et al. as they performed two usability studies narrowing the procedure to the certificate acquisition [5].

Finally, the study performed by Tiefenau et al. reveals that even experienced administrators struggle with keeping their systems up-to-date as the decision process that precedes the application of patches requires attention and is time-consuming [62]. To overcome this limitation, Li et al. suggests that helping system administrators during the information gathering would simplify the updating efforts and the likelihood of prioritizing the updates for the managed systems [38].

### 4 TOOL DESCRIPTION

With the goal of evaluating the benefits of providing mitigation hints in a TLS security report we need to identify an offline TLS analyzer to use for our user study. The market offers several tools to support sysadmins with security TLS configurations, all of them adopt a similar approach, i.e., they repeatedly connect to the target server and send specifically crafted ClientHello messages. By checking the server's responses (i.e. Server-Hello messages), these tools infer the server configuration and assess if it is affected by known vulnerabilities. However, their reports usually contain only the list of detected vulnerabilities and they offer little or no explanation on how to actually mitigate the detected weaknesses. Among the publicly available scanners we can mention ssl-enum-ciphers [32], ssllscan [49], TLSSLed [58], testssl.sh [65], SSL Server Test [48] and sslyze [18]. We chose testssl.sh as it is the most complete and covers the largest amount of required features (see Table 1).

#### 4.1 testssl.sh's Report

testssl.sh is a powerful open-source Bash script [65] that supports a wide range of state-of-the-art TLS-related checks. Checks include availability of ciphers and protocols, server preferences and an extensive set of information from the server certificate and its chain of trust. The report looks complete and quite verbose as it does not focus on the detected issues only, but it gives an overall

view of the server status including also the passed checks (see a fragment in Figure 1).

## 4.2 Actionable Reports

To understand if system administrators would benefit from a concise yet informative report (see Section 5) we took a subset of `testssl.sh` reports, removed all the passed and informative checks and added a number of descriptive elements collected by fetching information from both scientific literature and each vendor’s technical documentation. The resulting report provides a clear explanation of how the detected vulnerabilities can be exploited and a set of actionable mitigation measures that aim to thwart their impact, and operatively guide a system administrator in removing the found security defects. The provided mitigations are described at various levels of abstraction (see a fragment in Figure 2):

**Textual description:** natural language description of the TLS vulnerability and related mitigations (brief explanation of the actions to perform).

**Code snippet:** a fragment of code that can be copy-pasted into the webserver’s configuration to seamlessly fix the weakness. Together with the snippet, the report will provide a set of steps on how to find the correct file/line to edit.

## 5 DEFINITION OF THE EXPERIMENTAL FRAMEWORK

The *goal* of this study is to analyze the effect of providing a set of mitigations to system administrators with the *purpose* of evaluating the support offered by the actionable reports in patching a defective TLS configuration. The *quality focus* regards how mitigation hints increase the developer capability to correctly and quickly patch a defective TLS configuration. We thus formulate the following two research questions:

- RQ1. Do a textual description of the mitigation and the corresponding code snippet *increase* the likelihood of a *correct* patch to a defective TLS configuration by a system administrator?
- RQ2. Do a textual description of the mitigation and the corresponding code snippet *decrease* the time required by a system administrator to patch a defective TLS configuration?

The main perspective from which our experiment should be evaluated is how actionable information can enhance the identification and patching of insecure TLS configurations in terms of speed and precision. There are also other interesting point of views to consider such as those of (i) a researcher interested to empirically assess the benefit of hints in patching defective TLS configurations; or (ii) a project manager, who has to make a decision of which development/maintenance tools and procedures to adopt, in order to guarantee effective deployment of a correctly configured infrastructure.

The experimental settings have been designed following the template and guidelines by Wohlin et al. [7] to select participants and present their demographics (Section 5.1), to define the experimental design and select appropriate metrics and dependent/independent variables (Section 5.2), to identify the most appropriate statistical tests (Section 5.3) and, eventually, to identify the threats to the validity of our findings (Section 5.4).

## 5.1 Demographic Statistical Sample

We involved 62 participants in this study. They are Bachelor and Master students from the departments of Computer Science and Mathematics of the University of Trento playing the role of unexperienced system administrators who should patch defective TLS configuration files.

The study has been conducted as part of laboratory lectures in two courses of cybersecurity offered in our University.

Participants were aware that they could drop at any time with no consequences, as they would not have been evaluated for their performance in the experiment. There was no compensation (neither money nor bonus in the exam mark) for their participation in the study.

A profiling survey have been used to collect demographic data from the participants.

**Seniority.** The first question splits participants according to their seniority: 22 participants are Bachelor and 40 are Master students.

**Year.** 11 participants attend the 2nd and 11 the 3rd year of the Bachelor program, while 26 participants attend the 1st year and 14 the 2nd year of the Master program.

**Academic background.** We filled a list of the related University courses, whose content might have been relevant to influence the result of a corrective task on TLS configuration. The answers are shown in Table 2. Participants background was collected in terms of which related courses they already attended or not (column marked with  $\checkmark$  and  $\times$ , respectively, in the table). Most of the participants (i.e., 50 over 62) already attended the course *Introduction to Computer and Network Security*, while almost half of the participants (i.e., 34 out of 62) attended the course about *Security Testing*. Less participants attended *Cryptography* (21 students) and *Network Security* (10 students). A smaller group attended *Complexity, Crypto and Financial Technology* (3 participants), *Cyber-Security Risk Assessment* (5 participants) and *Offensive Security* (4 participants).

**Table 2: Demographics: Participants’ Academic background.**

Course	$\times$	$\checkmark$
Introduction to Computer and Network Security	12	50
Security Testing	28	34
Cryptography	41	21
Network Security	52	10
Cyber Security Risk Assessment	57	5
Offensive Security	58	4
Complexity, Crypto and Financial Technology	59	3

**Technical background.** Additionally, we collected the technical background of participants, in terms of which tasks they conducted in the past (data shown in Table 3). Only 15 are expert in manually configuring a TLS server because they already did it (column marked with  $\checkmark$ ), while half of them already configured Apache HTTP servers (33 participants) and created or edited other Unix configuration files (36 participants). The large majority of the participants are fluent in basic Unix administration tasks, such as navigating in the file system (61 participants), working with folders (60 participants), editing files (59 participants) and installing system packages (57 participants).

Testing vulnerabilities

```

Heartbleed (CVE-2014-0160)      not vulnerable (OK), timed out
CCS (CVE-2014-0224)           not vulnerable (OK)
Ticketbleed (CVE-2016-9244), experiment. not vulnerable (OK)
ROBOT                          Server does not support any cipher suites that use RSA key transport
Secure Renegotiation (CVE-2009-3555) not vulnerable (OK)
Secure Client-Initiated Renegotiation not vulnerable (OK)
CRIME, TLS (CVE-2012-4929)     VULNERABLE (NOT ok)
BREACH (CVE-2013-3587)        no HTTP compression (OK) - only supplied "/" tested
POODLE, SSL (CVE-2014-3566)    not vulnerable (OK)
TLS_FALLBACK_SCSV (RFC 7507)  No fallback possible, no protocol below TLS 1.2 offered (OK)
SWEET32 (CVE-2016-2183, CVE-2016-6329) not vulnerable (OK)
FREAK (CVE-2015-0204)         not vulnerable (OK)
DROWN (CVE-2016-0800, CVE-2016-0703) not vulnerable on this host and port (OK)
                                make sure you don't use this certificate elsewhere with SSLv2 enabled services
                                https://censys.io/ipv4?q=8EEC7DB209E2113BEC94EA55102C08B6CBCAA26317B1B73B7D41689DC6F93A66
LOGJAM (CVE-2015-4000), experimental not vulnerable (OK): no DH EXPORT ciphers, no DH key detected
BEAST (CVE-2011-3389)         no SSL3 or TLS1 (OK)
LUCKY13 (CVE-2013-0169), experimental potentially VULNERABLE, uses cipher block chaining (CBC) ciphers with TLS. Check patches
RC4 (CVE-2013-2566, CVE-2015-2808) no RC4 ciphers detected (OK)
    
```

Figure 1: testssl.sh report fragment

**CRIME**

By exploiting the information leakage provided by DEFLATE (compression algorithm), an attacker is able to retrieve the session cookie. In particular, the attacker guesses parts of the cookie, injects them in a valid client packet and analyzes the server's response. Thanks to the properties of a DEFLATE output, if the server's response is bigger than an untouched packet, then the guess is wrong.

**Mitigation**

Disable the TLS compression mechanism

**Code Snippet**

1. open your Apache configuration file (default: `/usr/local/apache2/conf/extra/httpd-ssl.conf`);
2. search for the line starting with: **SSLCompression**
  - if found, change the value to **off**;
  - if not, add the line `SSLCompression off`.

N.B. restart the server by typing: `sudo /usr/local/apache2/bin/apachectl -k restart`

Figure 2: Enhanced report fragment on CRIME

Table 3: Demographics: Participants' Technical background.

Technical skill	X	✓
Configure TLS servers	47	15
Configure Apache HTTP instances	29	33
Create/edit UNIX configuration files	26	36
Change working folder	1	61
Create/remove folder	2	60
Edit file	3	59
Install package	5	57

**Number of participants.** Establishing the right number of participants to a user study is a difficult task that can be completed only ex-post, after the experimental data are collected (by estimating the *power* of the test, as done in [11]). However, a replication with a large number of participants is mandatory for inconclusive studies, where the observed difference in independent variables was not statistically relevant (*power* was low). As a matter of fact, between 15 and 20 participants is considered reasonable to draw conclusions

from statistical analyses of the results [59]; in our case we have 4 times this number of participants.

Concerning the profiles of participants, we are aware that the expertise of students may be different from that of professionals. However, finding professionals available to conduct a demanding experiment as the one we designed is not easy. We mitigated this limitation by considering students with different levels of education (Bachelor and Master) and by making sure that participants had enough knowledge on TLS and its related vulnerabilities. All in all, the use of undergraduate students as a proxy of junior developers to draw conclusions is a common practice in empirical software engineering that is largely accepted and validated [28, 51, 61].

**Ethical considerations.** Participation was voluntary and students could have chosen not to attend the experiment without negative results on the final evaluation of the exam. Those students who opted to participate were aware that they would not be evaluated based on their performance and that they were receiving no compensation (neither money nor bonus in the exam mark) for the participation in the study; they were aware that they could

drop off at any time with no consequence. During all the experimentation, we strive to adhere to the general ethical principles stated in the ACM code of conduct [19] with particular attention to trustworthiness, fairness, privacy, and confidentiality.

## 5.2 Experimental Setup and Execution

**Systems.** The systems used to conduct the experiment are two web servers with defective TLS configurations, running Apache HTTP Server v2.4.37 and OpenSSL v1.0.2. Each incorrect configuration exposes the corresponding system to one specific attack. The two systems are:

- $S_1$ : a defective webserver vulnerable to BREACH; and
- $S_2$ : a defective webserver vulnerable to CRIME.

They are packaged as two distinct Virtual-box machines. Instead of using a random pair of detectable misconfigurations (e.g., POODLE [40], Sweet32 [6] or others) we selected two vulnerabilities that are comparable in terms of complexity of the operations required to patch. Moreover, we ensured that they can realistically be fixed in two hours, taking into account the student’s technical background. In particular, both are prone to the same type of information leakage caused by DEFLATE, but exploited using two different attacks (as discussed in Section 2.2). It is important to note that these systems are representative of realistic TLS configurations as both Apache and OpenSSL are respectively the most popular webserver [15] and TLS library [14]. To make the corrective tasks independent, only one vulnerability is present in each system.

**Metrics.** To measure the support of mitigation actions to conduct a corrective maintenance on TLS configurations (i.e., vulnerability detection and fix), we identified the following variables. The main factor of the experiment—that acts as an independent variable—is the presence of the mitigation hints during the execution of the task. In our experiment, the *base* treatment case  $TR_{list}$  consists of the bare list of vulnerabilities, as it is provided by the analysis tool `testssl.sh`; and  $TR_{hint}$  consists of the actionable reports, that include not only the list of vulnerabilities, but also a textual description of the mitigations and a code snippet to apply the mitigation.

Moreover, by adopting an approach similar to the one described in the ISO 9241 (Part 11) standard [29], we instrumented the experimental settings to measure the following dependent metrics:

- **Correctness** of each corrective task performed by participants, which corresponds to the *System Effectiveness* in [29] and thus examines the participants’ ability to complete a task. Participants could repeat the scans as many times as they like during the experimental session. However, to consider a task correct, the participants were supposed to run a final scan and to show the experimenter that the freshly generated report contains no vulnerability.
- **Time** taken to perform a corrective task on a defective TLS configuration, which corresponds to the *System Efficiency* in [29]. We collected such information by asking participants to fill in—while performing the experimental tasks—start and end time of each task.

Finally, the *System Satisfaction* in [29] to evaluate the overall usability of the TLS analyzer report is measured by a survey questionnaire (shown in Appendix A and analyzed in Section 6.3).

**Experimental Design.** We adopt a counter-balanced experimental design intended to fit two lab sessions of 40 minutes each. Participants are randomly assigned to four groups (despite they work alone) balanced based on their seniority, each one working in two labs on different systems with different treatments. The design allows for considering different combinations of *Systems* and *Treatments* in different order across *Labs* (see Table 4).

Table 4: Experimental design.

	Group A	Group B	Group C	Group D
Lab 1	$S_1 + TR_{hint}$	$S_2 + TR_{list}$	$S_2 + TR_{hint}$	$S_1 + TR_{list}$
Lab 2	$S_2 + TR_{list}$	$S_1 + TR_{hint}$	$S_1 + TR_{list}$	$S_2 + TR_{hint}$

**Experimental Procedure.** Before the experiment, participants were properly trained with lectures and exercises on TLS, to recall the required background [1]. The purpose of training is to make participants confident about the kind of tasks they are going to perform and the environment they will have available.

The experimental context has been set as similar as possible to a realistic scenario. As such, participants could run scanning tools as often as wanted, they could inspect the scan reports and browse the Internet to look for additional information. Moreover, by performing a dry-run test we ensured that the overall experiment can realistically be finished in two hours. The dry-run test helped us also to refine the survey questionnaires. Participants have been delivered the following material:

- two virtual machines:  $S_1$  and  $S_2$ ;
- two digital documents containing the instructions for each treatment (i.e.,  $TR_{list}$  and  $TR_{hint}$ );
- a printed page containing a recap of the lessons learned during the training phase (e.g., the commands used).

The experiment was carried out according to the following procedure. Participants had to:

- (1) Complete a pre-experiment profiling survey questionnaire;
- (2) For Lab 1: (i) mark the start time; (ii) perform the corrective task; (iii) mark the stop time;
- (3) Complete a survey questionnaire on the first lab;
- (4) For Lab 2: (i) mark the start time; (ii) perform the corrective task; and (iii) mark the stop time;
- (5) Complete a survey questionnaire divided in three parts: a *1st part* on the second lab, a *2nd part* on a comparison between the two labs, and a *3rd part* to collect feedback on  $TR_{hint}$  (the treatment with mitigation hints).

The pre-experiment profiling survey collects demographic data about the participants, such as their previous experience with Apache HTTP Server and their knowledge of the Bash command language; the complete survey is included in the replication package available online [1] and we have described the collected data in Section 5.1.

Each lab can be considered over, and, thus, a participant can mark the stop time, only after proving that the task was successfully completed or because the available time has expired.

We report the list of questions for the survey questionnaire in Appendix A and discuss the answers in Section 6.3. The survey questionnaires deal with cognitive effects of the treatments on

the behavior of the participants and perceived usefulness of the provided report.

### 5.3 Statistical Tests

We are interested to assess if the presence of mitigation hints has an impact on the *Correctness* and in the *Time* taken to fix defective TLS configurations. However, observed differences in the correctness of corrective tasks and time spent on them could be due to random variations or measurement errors. To test if the observed difference is statistically significant, we use sound statistical tests. As a common practice, we accept a 5% probability of committing type-I error, i.e., assessing that the difference is significant when it is actually due to random error. Practically, this setting defines the threshold  $\alpha = 0.05$ , for considering the result of a statistical test significant.

The lack of significance might mean also that there was an effect, but the effect was not observable (possibly because of a too small set of observations), rather than that the effect does not exist, i.e., we risk to commit a type-II error (nullification fallacy [35]). To quantify the probability of this problem, when the significance threshold is not reached, we can estimate the probability  $\pi$  of committing a type-II error as  $1 - \text{Power}$ , where *Power* is the statistical power of the adopted statistical test. As common practice, assume a threshold  $\beta = 0.20$  and we consider the power adequate when  $\pi < \beta$ .

The decision of which statistical tests to use was based on test applicability conditions and best practices recommended or commonly accepted in authoritative literature.

For each participant there are two distinct data points, one for the first lab and another one for the second lab, so we never perform multiple pairwise comparisons with overlapping data. Thus, there is not risk to inflate the family-wise error rate, and no correction factor (e.g., Bonferroni or Holm) is needed.

**Correctness.** To analyze the differences in terms of *Correctness*, we looked at the frequencies of correct/wrong tasks and we used a test on categorical data, because the tasks can be either correct (completed successfully) or incorrect (completed unsuccessfully). In particular, we used Fisher’s exact test [16] that is applicable on categorical data (correct/wrong answers). Fisher’s exact test is more accurate than the  $\chi^2$  test for small sample sizes, which is another possible alternative to test the presence of differences in categorical data. The same analysis was conducted by [9].

**Time.** To test the differences in *Time*, we perform the two-tailed Mann-Whitney U test on all samples [57]. This test is applicable to compare (time duration) samples of two populations. As a non-parametric test, Mann-Whitney U test does not require data to be normally distributed.

**Effect size.** To quantify the magnitude of differences among the two treatments, we used two kinds of effect size measures, the *odds ratio* for the categorical variable *Correctness* and the Cliff’s delta effect size [25] for *Time*. An odds ratio of 1 indicates that the condition or event under study is equally likely in both groups (participants using  $TR_{list}$  and those using  $TR_{hint}$ ). An odds ratio greater than 1 indicates that the condition or event is more likely in the first group. An odds ratio less than 1 indicates that the condition or event is less likely in the first group. For independent samples, Cliff’s delta provides an indication of the extent to which two

(ordered) data sets overlap, i.e., it is based on the same principles of the Mann-Whitney test. Cliff’s Delta ranges in the interval  $[-1, 1]$ . It is equal to +1 when all values of one group are higher than the values of the other group and  $-1$  when reverse is true. Two overlapping distributions would have a Cliff’s Delta equal to zero. The effect size is considered small for  $0.148 \leq d < 0.33$ , medium for  $0.33 \leq d < 0.474$  and large for  $d \geq 0.474$  [13].

**Co-factors.** The analysis of other factors (participants’ background, the system, the lab) that could have influenced the *Correctness* and *Time* is performed using the *Generalized Linear Mixed Model* [30] (GLMM for short). They extend the Generalized Linear Model (GLM) by adding random effects to the linear predictor (GLM only supports fixed effects). Random effects are particularly appropriate with repeated measures design, i.e., when different data points are collected for the same participant (in our design, each participant worked at two tasks). GLMM incorporates a number of different statistical models: ANOVA, ANCOVA, MANOVA, MANCOVA, ordinary linear regression, t-test and F-test. It consists in fitting a linear model of the *dependent* output variables (*Correctness* or *Time*) as a function of the *independent* input variables (all factors, including the treatment, i.e., the vulnerability detection tool). GLMM is capable of testing a dependent output variable (experiment outcome) on many input variables (factors) and it allows to test the statistical significance of the influence of each factor separately.

GLMM requires to specify the *exponential-family* distribution based on the domain of the outcome. We have chosen:

- the binomial family with logit link function to fit the *Correctness*, as it corresponds to a logistic regression that is appropriate for a binary outcome (correct/wrong task);
- a Gamma family for fitting the *Time*, as it is appropriate for fitting a time duration that can be a positive decimal value.

However, as other models could have been used to fit the *Time*, we use the *Akaike Information Criterion* (AIC for short) to check that the chosen model (i.e., the Gamma exponential-family distribution) was the most appropriate to fit our data [50]. AIC is founded on information theory and it entails balancing the trade-off between the goodness of fit of the model and the size of the model.

**Surveys.** Two statistical tests have been used on survey questionnaire. The Fisher’s exact test is used on categorical data, to compare the frequencies of yes/no answers for participants who worked with different tools. The Mann-Whitney U test is used to analyze answers to overall questions (not specific to any lab) when the answers were formulated using a Likert scale, checking for the null-hypothesis that the average answer was negative or neutral.

### 5.4 Threats to Validity

The main threats to the validity of this experiment belong to the internal, construct, conclusion and external validity threat categories. We discuss each one of them in the following.

*Internal validity* threats concern external factors that may affect the independent variable. The chosen design allowed us to control a number of factors, namely participants background, system and learning across experimental sessions. Participants were not aware of the experimental hypotheses, participants were not rewarded for

the participation in the experiment and they were not evaluated on their performance in doing the experiment.

*Construct validity* threats concern the relationship between theory and observation. They are mainly due to how we measure the correctness and duration of tasks. As described in Section 5.2, we considered real vulnerabilities and we used a sound procedure to objectively evaluate whether the fixes were correct.

The background of participants was estimated according to their academic background and their technical knowledge.

*Conclusion validity* threats concern the relationship between treatment and outcome. We used statistical tests to draw our conclusions on the correctness and the time required to fix TLS misconfigurations. The adopted statistical tests are particularly robust (i.e., they do not give false rejections of the null hypothesis) under deviations from normality.

*External validity* concerns the generalization of the findings. In our experiments we considered two major vulnerabilities related to TLS configuration, namely BREACH and CRIME. Although different vulnerabilities might occur, the results obtained with these vulnerabilities already support well our interpretations.

Our experiment exploited one real-world web application running in a web server (i.e., Apache HTTP). Despite we consider that this web server is representative of other web servers (e.g., nginx), in principle different results could be obtained for different web servers.

The study was performed in an academic environment, which may differ substantially from an industrial setup. However, we mitigate this threat by using subjects with different background and different seniority, including Bachelor and Master students, some of which with experience with TLS, Apache HTTP and Unix. Moreover, we considered the seniority as a factor to detect any influence on the results.

## 6 EXPERIMENTAL RESULTS

This section presents the data collected during the experimental validation. After analyzing data with sound statistical tests, we formulate answers to the two research questions stated at the beginning of Section 5.

### 6.1 Analysis of Correctness

Table 5 shows the distributions of correct/wrong answers when using testssl.sh’s reports ( $TR_{list}$ ) and the actionable reports ( $TR_{hint}$ ).

Almost all the participants were able to correctly patch the vulnerability when they were provided mitigations and code snippets, only one participant was not. Conversely, when participants were provided only with the list of vulnerabilities,<sup>2</sup> just 40 were able to complete a correct vulnerability patch.

We apply Fisher’s test to check if the observed trend is statistically significant. The difference is statistically relevant (p-value < 0.001, and we recall that we assume significance when p-value <  $\alpha$ , with  $\alpha = 0.05$ ) with a *large* effect size (odds ratio = 30).

Table 6 reports the analysis of correctness with GLMM. The model takes into account not only the effect of the main treatment (i.e., availability of mitigation hints) but all the other factors that

<sup>2</sup>The number of participants who worked with  $TR_{list}$  does not sum to 62, because two participants attended only the first lab.

**Table 5: Correctness in fixing an incorrect TLS conf.**

	$TR_{list}$	$TR_{hint}$
Correct	40 (67%)	61 (98%)
Wrong	20 (33%)	1 (2%)

we considered in our experimental design, i.e., the *System*, the *Lab*, the profile of the participants (the *Year* attended and their *Seniority*). The random-effects term is the participant who took parts in the two labs.

Statistically significant cases are in boldface. Consistently with the Fisher’s test, we can observe that the availability of mitigation hints significantly influences the correctness of a task related to fixing a TLS configuration file, as  $\Pr(>|z|) = 0.0010$ . However, it is the only significant factor.

The probability of committing type-II error obtained from GLMM power analysis is 0.05, which is smaller than 0.20. So, we can claim that the missing significance is not due to insufficient experimental data points, but to missing causal correlation between independent and dependent variable.

*System* does not significantly influence the *Correctness* of tasks, so we can conclude that the two applications and the two corrective tasks were well-balanced and none was harder or easier to fix. Moreover, the *Lab* is not a significant factor, thus there is no evident learning effect between the two experimental sessions. This means that having performed a first lab does not improve the accuracy in the second lab and we only measure the difference actually due to the independent variable (i.e., the presence of mitigation hints).

Based on these results, we can answer the research question RQ1 (see beginning of Section 5) as follows:

*Providing a text description together with a code snippet of the mitigation increases the capability of a system administrator to patch the defect in the TLS configuration. In fact, we observed that participants deliver correct fixes in 98% of the cases when this additional information has been included in the security reports, while the rate of correct fixes is 67% when no mitigation is provided.*

**Table 6: Analysis of Correctness (GLMM)**

	Estimate	Std. Error	z value	$\Pr(> z )$
(Intercept)	15.26	14.49	1.05	0.2923
Treatment	-14.12	4.28	-3.30	<b>0.0010</b>
System	-0.72	2.34	-0.31	0.7596
Lab	0.43	2.34	0.18	0.8534
Seniority	-6.07	13.78	-0.44	0.6598
Year	3.66	7.07	0.52	0.6048

### 6.2 Analysis of Time

We now analyze the time taken to fix a TLS configuration. Thus, we only consider time information for those participants who correctly fixed TLS configurations, and we discard data for incomplete and wrong tasks.

Figure 3 shows to box-plot of the time (in minutes) taken to fix a TLS configuration. Descriptive statistics (number of data points, mean, median and standard deviation) are summarized in Table 7.



On average, when `testssl.sh` is used ( $TR_{list}$ ) fixing a wrong configuration takes 23 minutes. When using the actionable reports ( $TR_{hint}$ ) the amount of time, on average, is reduced to less than 8 minutes.

According to the result of Mann-Whitney test, this difference is statistically significant ( $p$ -value < 0.001) with a *large* effect size (Cliff’s Delta = 0.8819).

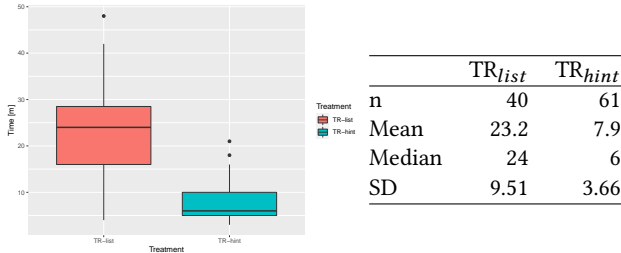


Figure 3 & Table 7: Time (in minutes) to fix a security issue.

Table 8 reports the analysis of *Time* with GLMM, with statistically significant cases in boldface. Consistently with the results of Mann-Whitney test, GLMM confirms that the availability of mitigation hints significantly influences the time needed to fixing a wrong TLS configuration file. Similarly to what previously observed for the *Correctness*, other factors have no significant influence on the *Time* to fix.

Table 8: Analysis of Time (GLMM)

	Estimate	Std. Error	t value	Pr(> z )
(Intercept)	0.04	0.03	1.39	0.1637
Treatment	0.09	0.01	12.55	<b>&lt;0.0001</b>
System	-0.01	0.01	-1.57	0.1176
Lab	0.01	0.01	1.02	0.3077
Seniority	-0.02	0.02	-0.79	0.4274
Year	0.00	0.01	0.39	0.6954

In this case, however, the probability of a type-II error obtained from GLMM power analysis is larger than 0.20. So, differently than the analysis of *Correctness*, the missing significance of cofactors influence (e.g., *System* and *Lab*) on *Time* cannot be interpreted as considerations on the experimental design.

Eventually, considering that this GLMM model could have been computed with different exponential-family distributions, we need to check that our choice was the most appropriate to fit our data. To this aim, we used the Akaike Information Criterion (AIC). This consists in fitting our data with other models and compare their AIC value. AIC values for other models are the following: 681.33 for Gaussian, 614.64 for Gamma and 631.43 for Poisson. As we can see, our choice is confirmed, because the AIC value for *Gamma* family is much lower than for the other models.

Considering these results, we can answer the research question RQ2 (see beginning of Section 5) as follows:

Providing a text description together with a code snippet of the mitigation decreases the time needed by a system administrator to patch the defect in the TLS configuration. In fact, we observed that in average it took 8 minutes to fix a misconfiguration when this additional information has been included in the security reports, while on average it took 23 minutes when no mitigation is provided.

Table 9: Analysis of survey questionnaire (Fisher’s test)

Question	$TR_{list}$		$TR_{hint}$		P-value
	Yes	No	Yes	No	
Enough time	40	20	61	1	<b>&lt;0.0001</b>
No difficulty	18	42	57	5	<b>&lt;0.0001</b>
Online search	53	7	3	59	<b>&lt;0.0001</b>
	<60%	≥60%	<60%	≥60%	
Configuration	53	7	41	21	<b>0.0048</b>
Documentation	24	36	61	1	<b>&lt;0.0001</b>

### 6.3 Analysis of Survey Questionnaire

The survey questionnaire (question formulation are reported in Appendix A) is composed of three parts and is meant to collect the participants opinion on the experiment.

*First Part.* The objective of the first part is to collect the participants opinion on the security reports, to indirectly compare them. Answers to the this first part are reported in Table 9. For the first three questions, the table reports the number of yes/no answers for participants who worked just with the list of vulnerabilities (2nd and 3rd columns, respectively) and with the mitigation actions within the actionable reports (4th and 5th columns).

The fourth and fifth questions asked participant to report the percentage of their lab time that they spent on specific tasks. Answers are in Likert scale (“< 20%”, “≥ 20% and < 40%”, “≥ 40% and < 60%”, “≥ 60% and < 80%” and “≥ 80%”). The table reports the number of participant who answered when a task took less than 60% or more than 60% of the lab time. The last column of Table 9 reports the significance (i.e.,  $p$ -value) computed by applying with the Fisher’s test for each question, to reveal statistical significance for the various answers. Significant cases are highlighted in boldface.

According to the first question, participants working with the actionable reports (i.e.,  $TR_{hint}$ ) considered that the time allocated to the task was enough, while time was short for participants assigned  $TR_{list}$ . This result is consistent with the analysis of *Time*, of Section 6.2, were participants who worked with no mitigation hints took longer to complete their tasks.

Consistently, looking at responses to the second question, we notice that only participants working with the actionable reports experienced no difficulty in completing the tasks. Tasks were harder to complete when participants were only supported by the list of security defects.

Considering the answers to the third question, we see that online searches were used by the majority of the participants who worked just with the list of vulnerabilities (53 positive answers versus 7 negative answers). When the actionable reports were used, instead, the majority of participants did not resort to online search (3 positive answers versus 59).

**Table 10: Analysis of survey questionnaire (Mann-Whitney test)**

Question	Strongly disagree	Disagree	Neutral	Agree	Strongly agree	P-value
Mitigations hints useful	0	3	15	21	20	< <b>0.0001</b>
Code snippets useful	1	1	10	15	32	< <b>0.0001</b>

Moving at the next two questions about time, we see a different approach on solving the assigned task. 53 participants assigned  $TR_{list}$  spent less than 60% of the lab time looking at TLS configuration code, meaning that they did not try to understand how TLS was configured but focused on searching online or in the TLS report for the solution. Conversely, when using the actionable reports almost none of the participants searched online for TLS documentation.

*Second Part.* The second part of the survey is a more direct comparison between the two alternative approaches. In fact, we asked participants to make an explicit decision about the two reports. For each question, Table 11 reports the number <sup>3</sup> of decisions that participants formulated about tasks supported by  $TR_{hint}$  and  $TR_{list}$ .

**Table 11: Analysis of survey questionnaire.**

Question	$TR_{list}$	$TR_{hint}$
Most useful	18	41
Most easy to read	7	52
Most complex to understand	53	6

The first question asked which was the most useful report when fixing the misconfiguration. The majority of the participants (41) considered the actionable reports the most useful.

The second question investigated if security reports were easy to read. The reports with mitigations were considered easier to read than the bare list of vulnerabilities.

The third question dealt with complexity of understanding. Consistently with the previous answers, the report that missed mitigations (i.e.,  $TR_{list}$ ) was considered more complex to understand than when mitigations were present (in  $TR_{hint}$ ).

Table 10 reports the answers to two other direct questions about mitigation hints and code snippets. Many participants strongly agree (20) or agree (21) that the textual mitigation hints were useful to complete the corrective task. A similar trend can be observed for the next question, participants strongly agree (32) or agree (15) that code snippet were useful to complete the corrective task. The result of Mann-Whitney test (null-hypothesis mean answer  $\leq$  “Neutral”) confirms that this trend is statistical significant.

*Third Part.* This last part, with only open questions, let participants write free text as feedback to the experiment. Its analysis would require a fundamentally different approach, mostly bases on *grounded theory* [21, 60] and is thus left for future work.

<sup>3</sup>We consider only 59 participants as one did not answer the 2nd part of the survey questionnaire.

## 7 LESSONS LEARNED

Often, new approaches may show trade-offs between contrasting goals and an optimal cost-benefit equilibrium has to be found. According to our experimental results, this is not the case when integrating actionable security hints into security reports. Indeed, we were able to observe a positive effect on both *Correctness* and *Time* of completion of the tasks.

In fact, we observed that mitigation hints help to patch defects in TLS configurations, by reducing the probability of error by 30 times and the time to complete the fix by 3 times.

In the following, we report the implications and general observations that we can formulate, based on the objective and quantitative results presented in the previous section.

**Limited information** *Automated security tools convey insufficient information.* A TLS configuration is quite complex as it contains a lot of properties that might not be of immediate understanding, thus the bare list of security problems is not informative enough to let a system administrator fix it. Indeed, additional information is needed to fill the gap of a security report and guide the system administrator towards figuring what changes are needed. In our experiment, the participants who received the list of vulnerabilities had to search online to understand how to fix security defects, while this was not required to those who worked with actionable hints (see Section 6.3 and Table 9).

**Correctness and Time** *Actionable maintenance hints improve correctness and time to fix.* Despite different tools find the same vulnerabilities, it is crucial how these are reported to system administrators. When actionable hints are available, a security report is more usable, user-friendly and able to guide system administrators towards a mitigation (see Table 6) in a timely manner (see Table 8 and Figure 3). Researchers and practitioners should keep this result in mind when developing new automated security tools. Scan results should be complemented with explanations and operational suggestions on how to solve the security problem or, at least, where to find additional information to guide towards the solution.

**Perceived effect** *Mitigations hints are easier to read and more useful than the list of vulnerabilities.* When conducting corrective maintenance, a flat list of the detected vulnerabilities is not perceived as very useful, probably because they are not very informative nor easy to read. This leaves system administrators with no clue where to look for getting the necessary information or to distinguish relevant from irrelevant data gathered from, e.g., online searches. Thus, additional information has to be collected either spending time in the code or reading additional documentation. Conversely, sysadmins are aware of the benefits of actionable mitigation hints because they are considered easier to read, more useful to support a fixing task and do not require additional time to fill knowledge gaps (see Survey Questionnaire and Table 11).

## 8 CONCLUSIONS

The usability of automated tools for vulnerability detection is quite a neglected topic. We have designed and conducted a user study aiming at filling this gap and answering the question in our title *Do Security Reports Meet Usability?* Our experiments reveal that the usability of reports is largely impacted by the availability of contextualized actionable hints, as they have a positive effect on both the correctness and the time needed to fix a TLS vulnerability. In addition, empirical evidence allows us to formulate a set of lessons learned that pave the way towards improving in general the usability of automated security tools. We can summarize the lesson learnt as follows. The main drawback is that the information provided (such as a flat list of detected problems) is often insufficient to enable users with little experience in security to mitigate a vulnerability. This can be alleviated by adding succinct textual explanations describing both the identified vulnerabilities and how they can be exploited in attacks. Besides reducing the time to fix a vulnerability and increasing the correctness of the applied patches, this approach has the potential to improve both knowledge and capabilities of administrators with less experience and to simplify the maintenance of complex systems. Thus, both the productivity and the security posture of an organization are improved. We used these considerations to build an open-source tool named TLSAssistant [54]. Together with the actionable reports we also added a learning dimension based on the concept of attack trees [53].

As future work, we plan to further improve TLSAssistant by increasing the amount of supported webserver, building new analysis modules and automatizing the changes suggested by the code snippets. We also plan to perform more experiments and understand how the application of suggested mitigations can jeopardize the availability of legacy systems and how participants may behave when presented with patches that require more complex tasks.

## REPRODUCIBILITY

All experimental material, including the vulnerable webservers, slides and questionnaires, together with the experimental data and the assets used for the training are available in the replication package [1].

## ACKNOWLEDGMENTS

This work has been partially supported by IPZS (Poligrafico e Zecca dello Stato Italiano, the Italian Government Printing Office and Mint) within the joint laboratory DigiMat Lab, the newly founded Futuro & Conoscenza S.r.l. and by the MIUR “Dipartimenti di Eccellenza” 2018-2022 grant.

## REFERENCES

- [1] 2021. Replication Package: Do Security Reports Meet Usability? Lessons Learned from Using Actionable Mitigations for Patching TLS Misconfigurations. <https://st.fbk.eu/complementary/ETACS2021>.
- [2] Yasemin Acar, Michael Backes, Sascha Fahl, Doowon Kim, Michelle L. Mazurek, and Christian Stransky. 2016. You Get Where You’re Looking for: The Impact of Information Sources on Code Security. In *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, San Jose, 289–305.
- [3] Luca Allodi, Silvio Biagioni, Bruno Crispo, Katsiaryna Labunets, Fabio Massacci, and Wagner Santos. 2017. Estimating the Assessment Difficulty of CVSS Environmental Metrics: An Experiment. In *Future Data and Security Engineering*, Tran Khanh Dang, Roland Wagner, Josef Küng, Nam Thoai, Makoto Takizawa, and Erich J. Neuhold (Eds.). Springer International Publishing, Cham, 23–39.
- [4] Luca Allodi, Marco Cremonini, Fabio Massacci, and Woo Hyun Shim. 2020. Measuring the accuracy of software vulnerability assessments: experiments with students and professionals. *Empirical Software Engineering* 25 (01 2020). <https://doi.org/10.1007/s10664-019-09797-4>
- [5] Matthew Bernhard, Jonathan Sharman, Claudia Ziegler Acemyan, Philip Kortum, Dan S. Wallach, and J. Alex Halderman. 2019. On the Usability of HTTPS Deployment <https://doi.org/10.1145/3290605.3300540>. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI ’19). Association for Computing Machinery, New York, NY, USA, 1–10. <https://doi.org/10.1145/3290605.3300540>
- [6] Karthikeyan Bhargavan and Gaëtan Leurent. 2016. On the Practical (In-)Security of 64-bit Block Ciphers: Collision Attacks on HTTP over TLS and OpenVPN. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24–28, 2016*. <https://doi.org/10.1145/2976749.2978423>
- [7] Michelle Cartwright. 2001. Book Review: Experimentation in Software Engineering: An Introduction. By Claes Wohlin, Per Runeson, Martin Höst, Magnus C. Ohlsson, Björn Regnell and Anders Wesslén. Kluwer Academic Publishers, 1999, ISBN 0-7923-8682-5. *Software Testing, Verification and Reliability* 11, 3 (2001), 198–199. <https://doi.org/10.1002/stvr.230>
- [8] Mariano Ceccato, Massimiliano Di Penta, Paolo Falcarin, Filippo Ricca, Marco Torchiano, and Paolo Tonella. 2014. A family of experiments to assess the effectiveness and efficiency of source code obfuscation techniques. *Empirical Software Engineering* 19, 4 (2014), 1040–1074.
- [9] Mariano Ceccato, Massimiliano Di Penta, Paolo Falcarin, Filippo Ricca, Marco Torchiano, and Paolo Tonella. 2014. A family of experiments to assess the effectiveness and efficiency of source code obfuscation techniques. *Empirical Software Engineering* 19, 4 (2014), 1040–1074.
- [10] Mariano Ceccato, Massimiliano Di Penta, Jasvir Nagra, Paolo Falcarin, Filippo Ricca, Marco Torchiano, and Paolo Tonella. 2009. The effectiveness of source code obfuscation: An experimental assessment. In *2009 IEEE 17th International Conference on Program Comprehension*. IEEE. <https://doi.org/10.1109/icpc.2009.5090041>
- [11] Mariano Ceccato, Alessandro Marchetto, Leonardo Mariani, Cu D Nguyen, and Paolo Tonella. 2015. Do automatically generated test cases make debugging easier? an experimental assessment of debugging effectiveness and efficiency. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 25, 1 (2015), 1–38.
- [12] Mariano Ceccato and Riccardo Scandariato. 2016. Static analysis and penetration testing from the perspective of maintenance teams. In *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. Association for Computing Machinery, New York, NY, USA, 1–6.
- [13] J. Cohen. 1988. *Statistical power analysis for the behavioral sciences (2nd ed.)*. Lawrence Earlbaum Associates, Hillsdale, NJ.
- [14] Datanyze. 2021. OpenSSL Market Share and Competitor Report <https://www.datanyze.com/market-share/other-it-infrastructure-software>.
- [15] Datanyze. 2021. Web and Application Servers Market Share Report <https://www.datanyze.com/market-share/web-and-application-servers>.
- [16] Jay L. Devore. 2007. *Probability and Statistics for Engineering and the Sciences*. Duxbury Press; 7 edition.
- [17] Thomas Dierks and Eric Rescorla. 2008. The Transport Layer Security (TLS) Protocol Version 1.2 <http://www.rfc-editor.org/rfc/rfc5246.txt>. Internet Requests for Comments.
- [18] Alban Diquet. 2021. Github: slyze <https://github.com/nabla-c0d3/sslyze>.
- [19] Association for Computing Machinery. 2018. ACM Code of Ethics and Professional Conduct <https://www.acm.org/binaries/content/assets/about/acm-code-of-ethics-booklet.pdf>.
- [20] Martin Georgiev, Subodh Iyengar, Suman Jana, Rishita Anubhai, Dan Boneh, and Vitaly Shmatikov. 2012. The most dangerous code in the world. In *Proceedings of the 2012 ACM conference on Computer and communications security - CCS ’12*. ACM Press. <https://doi.org/10.1145/2382196.2382204>
- [21] B. G. Glaser and A. L. Strauss. 1967. *The Discovery of Grounded Theory*. Aldine, Chicago.
- [22] Y. Gluck, N. Harris, and A. Prado. 2012. BREACH: reviving the CRIME attack <http://breachattack.com/>.
- [23] Peter Leo Gorski, Luigi Lo Iacono, Dominik Wermke, Christian Stransky, Sebastian Möller, Yasemin Acar, and Sascha Fahl. 2018. Developers Deserve Security Warnings, Too: On the Effect of Integrated Security Advice on Cryptographic API Misuse <https://www.usenix.org/conference/soups2018/presentation/gorski>. In *Fourteenth Symposium on Usable Privacy and Security (SOUPS 2018)*. USENIX Association, Baltimore, MD, 265–281.
- [24] Matthew Green. 2011. A diversion: BEAST Attack on TLS/SSL Encryption <https://blog.cryptographyengineering.com/2011/09/21/brief-diversion-beast-attack-on-tlssl/>.
- [25] Robert J. Grissom and John J. Kim. 2005. *Effect sizes for research: A broad practical approach* (2nd edition ed.). Lawrence Earlbaum Associates.

- [26] Norman Hänsch, Andrea Schankin, Mykolai Protsenko, Felix Freiling, and Zinaida Benenson. 2018. Programming Experience Might Not Help in Comprehending Obfuscated Source Code Efficiently <https://www.usenix.org/conference/soups2018/presentation/hansch>. In *Fourteenth Symposium on Usable Privacy and Security (SOUPS 2018)*. USENIX Association, Baltimore, MD, 341–356.
- [27] Scott Hollenbeck. 2004. Transport Layer Security Protocol Compression Methods <http://www.rfc-editor.org/rfc/rfc3749.txt>. Internet Requests for Comments.
- [28] Martin Höst, Björn Regnell, and Claes Wohlin. 2000. Using students as subjects—a comparative study of students and professionals in lead-time impact assessment. *Empirical Software Engineering* 5, 3 (2000), 201–214.
- [29] ISO 9241-11 2018. ISO 9241. Ergonomics of human-system interaction — Part 11: Usability: Definitions and concepts.
- [30] Jiming Jiang. 2007. *Linear and generalized linear mixed models and their applications*. Springer Science & Business Media.
- [31] Stephen Farrell Kathleen Moriarty, CIS. 2021. Deprecating TLSv1.0 and TLSv1.1 <https://tools.ietf.org/html/rfc8996>.
- [32] Mak Kolybabi and Gabriel Lawrence. 2020. ssl-enum-ciphers <https://nmap.org/nse/doc/scripts/ssl-enum-ciphers.html>.
- [33] Katharina Krombholz, Karoline Busse, Katharina Pfeffer, Matthew Smith, and Emanuel von Zezschwitz. 2019. "If HTTPS Were Secure, I Wouldn't Need 2FA" - End User and Administrator Mental Models of HTTPS. In *2019 IEEE Symposium on Security and Privacy (SP)*. 246–263.
- [34] Katharina Krombholz, Wilfried Mayer, Martin Schmiedecker, and Edgar Weippl. 2017. "I Have No Idea What I'm Doing" - On the Usability of Deploying HTTPS <https://www.usenix.org/conference/usenixsecurity17/technical-sessions/presentation/krombholz>. In *26th USENIX Security Symposium (USENIX Security 17)*. USENIX Association, Vancouver, BC, 1339–1356.
- [35] Anton Kühberger, Astrid Fritz, Eva Lermer, and Thomas Scherndl. 2015. The significance fallacy in inferential statistics. *BMC research notes* 8, 1 (2015), 84.
- [36] Katsiaryna Labunets, Fabio Massacci, Federica Paci, Sabrina Marczak, and Flávio Moreira de Oliveira. 2017. Model comprehension for security risk assessment: an empirical comparison of tabular vs. graphical representations. *Empirical Software Engineering* 22, 6 (Feb 2017), 3017–3056. <https://doi.org/10.1007/s10664-017-9502-8>
- [37] Katsiaryna Labunets, Fabio Massacci, Federica Paci, and Le Minh Sang Tran. 2013. An Experimental Comparison of Two Risk-Based Security Methods. In *2013 ACM / IEEE International Symposium on Empirical Software Engineering and Measurement*. IEEE, 163–172. <https://doi.org/10.1109/esem.2013.29>
- [38] Frank Li, Lisa Rogers, Arunesh Mathur, Nathan Malkin, and Marshini Chetty. 2019. Keepers of the Machines: Examining How System Administrators Manage Software Updates For Multiple Machines. In *Fifteenth Symposium on Usable Privacy and Security (SOUPS 2019)*. USENIX Association, Santa Clara, CA. <https://www.usenix.org/conference/soups2019/presentation/li>
- [39] Microsoft-Inria. 2014. Triple Handshakes Considered Harmful: Breaking and Fixing Authentication over TLS <https://www.mitls.org/pages/attacks/3SHAKE>.
- [40] Bodo Möller, Thai Duong, and Krzysztof Kotowicz. 2014. This POODLE Bites: Exploiting The SSL 3.0 Fallback <https://www.openssl.org/~bodo/ssl-poodle.pdf>.
- [41] Mozilla Security. 2018. Web Security Cheat Sheet [https://infosec.mozilla.org/guidelines/web\\_security](https://infosec.mozilla.org/guidelines/web_security).
- [42] Alena Naiakshina, Anastasia Danilova, Eva Gerlitz, Emanuel von Zezschwitz, and Matthew Smith. 2019. "If You Want, I Can Store the Encrypted Password": A Password-Storage Field Study with Freelance Developers <https://doi.org/10.1145/3290605.3300370>. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300370>
- [43] Alena Naiakshina, Anastasia Danilova, Christian Tiefenau, Marco Herzog, Sergej Dechand, and Matthew Smith. 2017. Why Do Developers Get Password Storage Wrong? A Qualitative Usability Study <https://doi.org/10.1145/3133956.3134082>. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security* (Dallas, Texas, USA) (CCS '17). Association for Computing Machinery, New York, NY, USA, 311–328. <https://doi.org/10.1145/3133956.3134082>
- [44] NIST. 2012. CVE-2012-4929 <https://nvd.nist.gov/vuln/detail/CVE-2012-4929>.
- [45] NowSecure. 2017. Fully Validate SSL/TLS <https://books.nowsecure.com/secure-mobile-development/en/sensitive-data/fully-validate-ssl-tls.html>.
- [46] Juan C. Perez. 2016. SSL: Deceptively Simple, Yet Hard to Implement <https://blog.qualys.com/product-tech/2016/12/12/ssl-deceptively-simple-yet-hard-to-implement>.
- [47] Qualys. 2021. SSL Pulse <https://www.ssllabs.com/ssl-pulse/>.
- [48] Qualys. 2021. SSL Server Test <https://www.ssllabs.com/sslttest/>.
- [49] rbsec. 2017. sslscan <https://github.com/rbsec/sslscan/releases/tag/1.11.11-rbsec>.
- [50] Benjamin Saefken, Thomas Kneib, Clara-Sophie van Waveren, Sonja Greven, et al. 2014. A unifying approach to the estimation of the conditional Akaike information in generalized linear mixed models. *Electronic Journal of Statistics* 8, 1 (2014), 201–225.
- [51] İlaah Salman, Ayse Tosun Misirli, and Natalia Juristo. 2015. Are Students Representatives of Professionals in Software Engineering Experiments? (*ICSE '15*). IEEE Press, Florence, Italy, 666–676.
- [52] Riccardo Scandariato, James Walden, and Wouter Joosen. 2013. Static analysis versus penetration testing: A controlled experiment. In *2013 IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE. <https://doi.org/10.1109/issre.2013.6698898>
- [53] Bruce Schneier. 1999. Attack Trees [https://www.schneier.com/academic/archives/1999/12/attack\\_trees.html](https://www.schneier.com/academic/archives/1999/12/attack_trees.html).
- [54] Security & Trust Research Unit. [n.d.]. TLSAssistant <https://github.com/stfbk/tlsassistant>.
- [55] Amazon Web Services. 2021. Alexa Top Sites <https://aws.amazon.com/alexa-top-sites/>.
- [56] Yaron Sheffer, Ralph Holz, and Peter Saint-Andre. 2015. Summarizing Known Attacks on Transport Layer Security (TLS) and Datagram TLS (DTLS) <http://www.rfc-editor.org/rfc/rfc7457.txt>. Internet Requests for Comments.
- [57] David J. Sheskin. 2007. *Handbook of Parametric and Nonparametric Statistical Procedures (4th Ed.)*. Chapman & All.
- [58] Raul Siles. 2013. TLSSLed v1.3 <http://blog.taddong.com/2013/02/tlssled-v13.html>.
- [59] Janet M. Six and Ritch Macefield. 2016. How to determine the right number of participants for usability studies <https://www.uxmatters.com/mt/archives/2016/01/how-to-determine-the-right-number-of-participants-for-usability-studies.php>.
- [60] A. Strauss and J. Corbin. 1990. *Basics of Qualitative Research: Grounded Theory Procedures and Techniques*. Sage, London.
- [61] Mikael Svahnberg, Aybüke Aarum, and Claes Wohlin. 2008. Using Students as Subjects - an Empirical Evaluation <https://doi.org/10.1145/1414004.1414055>. In *Proceedings of the Second ACM-IEEE International Symposium on Empirical Software Engineering and Measurement (Kaiserslautern, Germany) (ESEM '08)*. Association for Computing Machinery, New York, NY, USA, 288–290. <https://doi.org/10.1145/1414004.1414055>
- [62] Christian Tiefenau, Maximilian Häring, Katharina Krombholz, and Emanuel von Zezschwitz. 2020. Security, Availability, and Multiple Information Sources: Exploring Update Behavior of System Administrators. In *Sixteenth Symposium on Usable Privacy and Security (SOUPS 2020)*. USENIX Association, 239–258. <https://www.usenix.org/conference/soups2020/presentation/tiefenau>
- [63] Alessio Viticchié, Leonardo Regano, Cataldo Basile, Marco Torchiano, Mariano Ceccato, and Paolo Tonella. 2020. Empirical assessment of the effort needed to attack programs protected with client/server code splitting. *Empirical Software Engineering* 25, 1 (2020), 1–48.
- [64] Alessio Viticchié, Leonardo Regano, Marco Torchiano, Cataldo Basile, Mariano Ceccato, Paolo Tonella, and Roberto Tiella. 2016. Assessment of Source Code Obfuscation Techniques. In *2016 IEEE 16th International Working Conference on Source Code Analysis and Manipulation (SCAM)*. IEEE, Los Alamitos, CA, USA, 11–20. <https://doi.org/10.1109/scam.2016.17>
- [65] Dirk Wetter. 2021. /bin/bash based SSL/TLS tester: testssl.sh <https://testssl.sh>.

## A QUESTIONNAIRES' CONTENT

Table 12 shows the three parts of the survey questionnaire.

**Table 12: Questions in the survey questionnaire.**

#	Question	Answer
<i>First part</i>		
<i>Lab 1</i>		
1	I had enough time to perform the task	[1-5] Likert scale
2	I experienced no difficulty in patching the vulnerability given the report	[1-5] Likert scale
3	How much time (in terms of percentage) did you spend looking at the TLS configuration code	[ $\leq a$ ] Likert scale
4	How much time (in terms of percentage) did you spend looking at online documentation on TLS vulnerabilities	[ $\leq a$ ] Likert scale
5	Provide some examples of online queries you used to search the vulnerabilities online (e.g. keywords used)	[None]/free text
6	Which steps did you take to perform the tasks? (e.g. run command Y, opened file X, ..)	[None]/free text
<i>Lab 2</i>		
1	I had enough time to perform the task	[1-5] Likert scale
2	I experienced no difficulty in patching the vulnerability given the report	[1-5] Likert scale
3	How much time (in terms of percentage) did you spend looking at the TLS configuration code	[ $\leq a$ ] Likert scale
4	How much time (in terms of percentage) did you spend looking at online documentation on TLS vulnerabilities	[ $\leq a$ ] Likert scale
5	Provide some examples of online queries you used to search the vulnerabilities online (e.g. keywords used)	[None]/free text
6	Which steps did you take to perform the tasks? (e.g. run command Y, opened file X, ..)	[None]/free text
<i>Second part</i>		
1	Which report did you find more useful?	Lab1/Lab2
2	Which report did you find more easy to read?	Lab1/Lab2
3	Which report did you find more complex to understand?	Lab1/Lab2
4	The textual description of the mitigation is useful to complete the tasks	[1-5] Likert scale
5	The code snippet is useful to complete the tasks	[1-5] Likert scale
6	How did you use the code snippet?	[None]/free text
<i>Third part</i>		
1	Would you use it for your work?	[None]/free text
2	Motivate your answer (to the previous question)	[None]/free text
3	Do you know any tool that performs similar tasks?	[None]/free text
4	Do you have any suggestion related to the tool usage?	[None]/free text
5	Do you have any suggestion related to the amount of information provided by the tool's report (Report.md)?	[None]/free text