

Human Aspect of Threat Analysis: A Replication

Katja Tuma

k.tuma@vu.nl

Vrije Universiteit Amsterdam

The Netherlands

Winnie Mbaka

w.mbaka@vu.nl

Vrije Universiteit Amsterdam

The Netherlands

ABSTRACT

Background: Organizations are experiencing an increasing demand for security-by-design activities (e.g., STRIDE analyses) which require a high manual effort. This situation is worsened by the current *lack of diverse (and sufficient)* security workforce and inconclusive results from past studies. To date, the deciding human factors (e.g., diversity dimensions) that play a role in threat analysis have not been sufficiently explored.

Objective: To address this issue, we plan to conduct a series of exploratory controlled experiments. The main objective is to empirically measure the human-aspects that play a role in threat analysis alongside the more well-known measures of analysis performance.

Method: We design the experiments as a differentiated replication of past experiments with STRIDE. The replication design is aimed at capturing some similar measures (e.g., of outcome quality) and additional measures (e.g., diversity dimensions). We plan to conduct the experiments in an academic setting.

Limitations: Obtaining a balanced population (e.g., wrt gender) in advanced computer science courses is not realistic. The experiments we plan to conduct with MSc level students will certainly suffer this limitation.

KEYWORDS

Threat Analysis, Human Aspects, Empirical Software Engineering, Replication, Controlled Experiment

ACM Reference Format:

Katja Tuma and Winnie Mbaka. 2022. Human Aspect of Threat Analysis: A Replication. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Security-by-design techniques [8, 30] have been used to prevent costly security fixes to software in later stages of the development life-cycle by analyzing security already during the design phase. Practitioners use threat analysis [36] to look for potential security threats in their product's software architecture. For instance, STRIDE [33] is a popular technique developed by Microsoft.

There is an increasing need to perform such architectural security analyses (e.g., latest BSIMM study reports an increased investment by more than 65% [10]) as the threat landscape evolves. However, threat analysis requires a high manual effort [31], demands the involvement of security and domain experts [6], and has been proven time and again difficult to fully automate [39].

Threat analysis practices are set back by a globally recorded shortage of the security workforce [4, 7]. In addition, the current security workforce is not diverse (e.g., with respect to gender) which may be viewed as an opportunity for a change.

Risk decisions (which are core to threat analysis) are made in face of uncertainty [3], thus there is space for subjective (and possibly biased) judgement [5, 15]. Empirical evidence of threat analysis performance indicators is a crucial piece of the puzzle to improve the situation. But, past empirical studies were either inconclusive about some performance indicators [38] or have focused on measuring performance indicators irrespective of the human factors [31, 37, 40]. Yet measuring such human factors is pivotal to understanding how to close the security workforce gap in the future.

To address these issues, we plan to conduct a series of exploratory controlled experiments with the aim of empirically measuring the human-aspects that play a role in threat analysis. In particular, we design a differentiated replication [21], where we capture some similar measures used in previous experiments [38] but also different measures (e.g., participant gender, nationality, type of outcomes).

2 RELATED WORK

We positioned our contributions with respect to existing literature on empirical studies of STRIDE and related replication studies.

Empirical studies of threat analysis. Several works have investigated STRIDE empirically. Scandariato et al. [31] performed a descriptive analysis to quantify the cost and effectiveness of STRIDE by measuring the productivity, precision, and recall in an academic setting. Their study reports students having higher rates in precision and recall with lower productivity rates.

Two studies [2, 6] conducted case studies investigating the challenges of STRIDE. Bernsmed et al. [6] performed an exploratory case study with the goal of investigating the challenges facing adoption of threat modeling using the Microsoft approach with STRIDE. The study was done in a company comprising five agile development projects. Their analysis elicited 21 challenges to threat modelling which were then mapped to existing literature and concluded that proper understanding of threat elicitation is required in order to actualise the functions of STRIDE especially in agile development.

Stevens et al. [34] conducted a qualitative case study to evaluate the impact of introducing threat modeling to an organization that

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.

Conference'17, July 2017, Washington, DC, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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

had previously not used it by applying the Center of Gravity (CoG) framework. The CoG is a risk-first threat analysis technique that has not been extensively used to analyze software security. The authors goal was to measure effectiveness and efficiency of CoG, using surveys and classroom sessions that involved 25 practitioners. The authors reported a very high accuracy from the participant's results, similar to other studies that have been conducted with experts. Different attack and risk-centric threat analysis techniques (such as Attack trees, CORAS, MUCs, to name a few) have also been investigated with empirical rigour. We point the interested reader to a systematic review for more details [36].

Replications. We frame our plan as a series of experimental replications. Generally, the goal of replication studies is to validate the experimental procedures of the original study using a different participation pool. This studies are aimed at generating new data [29] (as opposed to re-analyzing the same data in reproduction studies). We briefly mention some related replication studies.

In their original study, Labunets et al [20], compared two risk assessment methods, a visual and a textual method and reported that the visual method was more effective for identifying threats than the textual one. The same study was replicated in [19], applying similar procedures. In contrast to the original study, the replication reported that the two methods being investigated were (statistically) equivalent with regards to the quality of identified threats and security controls.

Several studies have empirically compared [17, 18, 23] and conducted replications [16, 24, 26] requirement engineering techniques (e.g., requirements elicitation). For brevity, we direct the interested reader a comprehensive review by Ambreen et al. [1].

3 RESEARCH QUESTIONS

Due to the academic setting we limit this study on observing gender, background, and nationality diversity dimensions (and exclude seniority). The main goal of this study is to measure the existence (or absence) of diversity effects on the actual and perceived analysis outcomes. Accordingly, we developed two research questions and hypotheses about each measure.

RQ1. *What is the effect of gender, background, and nationality on the **actual** threat analysis outcomes?*

To investigate RQ1, we pose hypotheses about the **equivalence of the sample means for the analysis outcomes**.

$$H1_1 : \text{Comp.Sci}_F = \text{Comp.Sci}_M$$

Regarding gender, we expect that the outcomes reported by women are equivalent to the outcomes reported by men. Studies of risk perception suggest that women perceive certain risks differently compared to men. Though we do not foresee strong differences, we might find some effects when it comes to risk priority.

$$H1_2 : \text{Comp.Sci}_1 = \text{Comp.Sci}_2$$

Regarding education, we expect that the students of various specialization tracks report equivalent outcomes for the same system under analysis.

$$H1_3 : \text{Comp.Sci}_{Na} = \text{Comp.Sci}_{Nb}$$

We expect that the students of various race and nationality report statistically equivalent outcomes.

RQ2. *What is the effect of gender, background, and nationality on the **perceived** threat analysis outcomes?*

To investigate RQ2, we pose hypotheses about the **equivalence of the sample means for the perceived analysis outcomes**.

$$H2_1 : \text{Perc}(\text{Comp.Sci}_F) < \text{Perc}(\text{Comp.Sci}_M)$$

Due to low confidence levels of female computer science students, we expect that the perceived quality of outcomes reported by women is lesser compared to the perceived quality of outcomes reported by men (regardless of the actual outcomes by both groups).

$$H2_2 : \text{Perc}(\text{Comp.Sci}_1) = \text{Perc}(\text{Comp.Sci}_2)$$

Regarding education, we expect that overall the students of various specialization tracks do not differ in their perceived quality of the outcomes they produced. We may find higher confidence levels of perceived quality for students that are following a security specialization track.

$$H2_3 : \text{Perc}(\text{Comp.Sci}_{Na}) = \text{Perc}(\text{Comp.Sci}_{Nb})$$

We expect that the students of various nationality do not differ in their perceived quality of the outcomes they produced.

4 REPLICATION PROTOCOL

4.1 Variables

Table 1 shows the variables of the study.

4.1.1 Independent. Gender is an individual's identity based on their sex, which is typically, man and woman, but can also be non-binary. Several studies [25] [13] [41] found evidence of bias against women in some software engineering communities, and sometimes negative perceptions about women working in teams. Thus, this is an interesting dimension to further investigate in the context of security.

Education is an individuals' level of achievement within a specific area of specialization (e.g., computer security vs AI) in academic studies. Typically, risk-based decisions in organizations have been made by persons in managerial positions who are assumed to have a better understanding of the product. However, technical skills of security experts or engineers should be taken into consideration during critical decision-making. Therefore, it is interesting to investigate the effect of education by including participants from different academic backgrounds (e.g., in communication sciences). Although in [14] study, education was found to be a non-significant variable, it is not clear whether this dimension has an impact in performing a RA task.

Nationality is generally used to refer to someone's country of origin, however, it may be coupled with other identifying aspects, such as the culture and language. On the other hand, race is a social construct that is associated with an individuals physical appearance, such as their skin colour. Determining the effect of nationality bias in security practices is to date an open question. In a previous study, Thomas et al. [35] conducted semi-structured interviews with 11 Black women in computing and report that Black women experienced a number of challenges, such as discrimination, expectations from others that are too high or too low, isolation, sexism, and racism. However, it is not clear whether some of

these challenges were as a result of their gender, race or nationality). But, few studies have focused on nationality diversity in the software engineering discipline [25].

4.1.2 Dependent. Since the quality of analysis lacks a formalised definition (e.g., often natural language is used to describe attack scenarios and informal notations are used for modeling [36]), we will use measures that can be easily reproduced. Namely, we can observe how diversity dimensions effect the *type of analysis outcomes*. Table 1 (dependant variables) shows various outcomes types that we observe.

Threats. We use the STRIDE threat categories to distinguish different type of threats. Tuma et al. [37, 38] noticed that expert analysis tend to be more balanced in terms of their review of different threat categories, while non-experts tend to report a high number of tampering, denial of service and information disclosure threats. We are interested to observe whether other diversity dimensions have an effect on category distribution of identified threats.

Assumptions. Assumptions are statements about the domain that may or may not be true. Assumptions are often implicit and dynamic in nature (i.e., they can be invalidated and modified as the project evolves). Van Landuyt and Joosen [40] find that the majority of assumptions (created by students during STRIDE) were used to either justify an existence of threats or are used to eliminate threats. In [40] a substantial subset (78%) of the assumptions was in direct reference to security-related concepts (i.e., security assumptions), however also domain assumptions (statements about component functionalities) were made. Thus, we are interested to investigate the effect of diversity dimensions on the type of assumptions.

Attack surface. defining an attacker profile and the *attack surface* are essential in determining the feasibility of an attack scenario. To this end, security analysis are expected to make these distinction prior to performing a threat analysis.

Risk priority. We refer to risk as a product of threat probability and impact. How individuals assess risk priorities may be related to their risk perception which is already well understood [11]. Since the number of identified threats explodes in realistic projects, practitioners must choose which threats are most urgent to mitigate. Thus they prioritize them based on estimations of risk.

Mitigations. There are different approaches that can be applied while mitigating security risks. Preventative (e.g., implementation of two-factor authentication), detective/reactive (e.g., using intrusion detection and access revocation techniques) and corrective (e.g., maintaining audit trails or restoring from a secure state). Security analysts can implement multiple strategies depending on domain-related factors, such as the cost associated with a specific mitigation strategy. Some diversity dimension (e.g., gender) may underestimate the ease with which a mitigation is actually implemented, as observed in [42]. Thus, it is interesting to observe how other diversity dimensions, including gender effect the type of mitigations that are identified during threat analysis.

4.2 Material

Training. In the first part of the training the participants will be introduced to some key security topics (such as CIAA triad, security threats, attack surface and vulnerabilities, security controls

and risk mitigations). The second part of the training will prepare the students to actually perform a threat analysis using one of the technique variants. The third part of the training will introduce the participants to the case study which will be the object of their analysis.

Case study documentation. We will use the same case study as in the original study. The home monitoring system (HomeSys) is an automated surveillance system designed for residential places. Its main objective is to enable the home-owner to remotely monitor their property. A detailed documentation of the case (requirements, architectural design, etc) will be made available to the participants.

Ground truth analysis. We will use one 'golden standard' data flow diagram and its' corresponding ground truth STRIDE analysis of the HomeSys case study from [38]. Since we do not aim to measure the quality of the diagrams created, and the DFD building is less time consuming compared to threat identification, we will provide a model to the participants. This will significantly simplify the comparison of the identified security threats. Similarly, we will provide the ground truth analysis to the participants that will be prioritizing threats and identifying security mitigations.

4.3 Task

The participants will be asked to individually fill-in a survey. The survey consists of three parts. First, a few questions about the students gender, background, nationality. Half of the participants will be asked to perform a STRIDE analysis (i.e., identify security threats). In contrast to previous studies, our participants will *be given the same graphical model* of HomeSys to analyze and they will analyze it using the same STRIDE technique. The other half of the participants will be asked to prioritize a list of security threats and identify security mitigations to high-priority threats. In contrast to previous studies, our participants will be given the graphical model *and the list of security threats*. To guide the threat identification the participants will use the documentation of STRIDE. Similar to the past studies, we will hand out a threat template csv to standardize the format of the outcomes reported. The participants will submit the files using the same survey. Finally, they will be asked a few questions regarding their perception of the task. Time taken to complete the task was captured using an online survey tool.

4.4 Participants

Our population is computer science students, with some differences in the elective courses and program choices (e.g., we plan to include students from various master programs, such as IA, computer Security, and Software Engineering). All participants are students enrolled in a course taught by the experimenters. At the beginning of the course we plan to hand out an entry survey to measure participants' background and areas of expertise relevant to the study. We expect most to be new to secure design techniques (e.g STRIDE, threat modeling, Data Flow Diagrams, misuse cases, attack trees etc). In addition, we expect the participants are unfamiliar with architectural modeling techniques (e.g sequence, component and deployment diagrams).

Name	Description	Scale	Operationalization
<i>Independent variables (design)</i>			
Gender	obtained from the gender of participants	nominal	multiple choice
Background	the program specialization and extra curriculum activities	nominal	multiple choice
Nationality	obtained from the nationality of participants	nominal	multiple choice
<i>Dependent variables</i>			
<i>**Different measures compared to existing literature**</i>			
Type of identified threats	distribution of categories of threats (spoofing, tampering, information disclosure, denial of service, elevation of privilege) that have been identified by the participants	nominal	see Section 4.2
Type of assumptions	distribution type of assumptions (domain, security) that have been reported by the participants	nominal	see Section 4.2
Type of attacks surface	distribution attack surfaces (physical, close-proximity, remote) of the identified threats	nominal	see Section 4.2
Risk priorities	distribution of risk priorities (high, medium, low) assigned to identified threats	nominal	see Section 4.2
Type of mitigations	distribution of type of identified mitigations (preventative, detective/reactive, corrective)	nominal	see Section 4.2
<i>Treated/Measured variables</i>			
Time spent on task	time (in hours) to complete the task using the prescribed technique	ordinal	automatically measured by the submission tool
Perceived precision (PP)	self-reported ratio between the number of correctly identified threats and all <i>threats identified</i>	ordinal	5-point Likert scale
Perceived recall (PR)	self-reported ratio between the number of correctly identified threats and all <i>existing</i> threats identified	ordinal	5-point Likert scale
Perceived usefulness (PU)	self-reported usefulness of the prescribed technique	ordinal	5-point Likert scale
Experience with security and modeling	self-reported experience in number of years or previously completed courses	ordinal	5-point Likert scale
Experience with STRIDE	self-reported experience in number of years or previously completed courses	ordinal	5-point Likert scale
Experience with domain of application	self-reported experience in number of years or previously completed courses	ordinal	5-point Likert scale
<i>**Different measures compared to existing literature**</i>			
Perceived cognitive load	the reported cognitive load (complexity) of the task using the prescribed technique	ordinal	5-point Likert scale
Perceived quality and efficacy	the reported quality and self efficacy of work conducted	ordinal	5-point Likert scale

Table 1: Variables of the differentiated replication experiments

4.5 Execution plan

Work division. The participants will be randomly divided into two groups (A and B). Group A will be tasked with analysing a provided data flow diagram of the HomeSys case study using STRIDE. Group B will be tasked with prioritizing a provided list of security threats and identifying security mitigations for high-priority threats. These treatment groups are formed only to divide the work to avoid overloading the participants performing an overly complex task individually.

Training. The participants will undergo an obligatory training lectures (about 3 hours) covering the topics mentioned above.

Hand-outs. After the training, participants will be given digital copies of all the support material (inc. lecture slides, case documentation, technique documentation, etc).

Physical labs. The experiment will be conducted during a four hour physical lab. The participants will be separated into different classrooms depending on their treatment group (to avoid spillover effects). Each classroom will be supervised by either a teaching assistant or the experimenters. Only questions about the experiment protocol will be answered.

Reports. The data will be collected through an online survey tool.

4.6 Analysis plan

Data cleaning. We will perform a preliminary check of the collected data. This will include removing submissions for which we did not get explicit consent by the participant. Second, we will remove clearly insincere submission attempts (if any).

TOST analysis of equivalence. We will use both difference and equivalence statistical tests. It is likely that we do not obtain normally distributed samples, thus we plan to use a non-parametric, Mann-Whitney test. TOST was initially proposed by [32] and is widely used in pharmacological and food sciences to check whether two treatments are equivalent within a specified range δ [9, 22]. Wherever possible (e.g., for Likert-scale questions) we will define the delta empirically. For instance by pooled variance σ_p across several samples reported in the literature on security risk analysis (e.g., in a four year interval) on variables ranging over a 5-item Likert scale for demographic statistics as to account for natural variability of the data.

Validity threats. There is typically around 20% (or less) female students enrolled in computer science programs. We are aware of the validity threats caused by an unbalanced population sample, which is omnipresent in all gender diversity studies in STEM disciplines [25]. To partially mitigate this threat, we will rally female computer scientist students towards participation through local feminist groups and similar community organized channels.

Since we do not include practitioners in this study, we can not observe the full complexity of the diversity effects (e.g., including seniority) that are actually present in organizations where threat analysis is routinely performed. Still, studies have shown [12, 27, 28] that the differences between the performance of professionals and graduate students are often limited.

We considered the threat of overloading the participants with a complex task. We mitigate this threat by splitting the participants into two groups, so individual participants get to either only focus on finding threats or focus on mitigating risks.

5 ACKNOWLEDGMENTS

REFERENCES

- [1] Talat Ambreen, Naveed Ikram, Muhammad Usman, and Mahmood Niazi. 2018. Empirical research in requirements engineering: trends and opportunities. *Requirements Engineering* 23, 1 (2018), 63–95.
- [2] Karin Bernsmed and Martin Gilje Jaatun. 2019. Threat modelling and agile software development: Identified practice in four Norwegian organisations. In *2019 International Conference on Cyber Security and Protection of Digital Services (Cyber Security)*. IEEE, 1–8.
- [3] Vicki Bier. 2020. The Role of Decision Analysis in Risk Analysis: A Retrospective. *Risk Analysis* 40, S1 (2020), 2207–2217.
- [4] Borka Jerman Blažič. 2021. Cybersecurity Skills in EU: New Educational Concept for Closing the Missing Workforce Gap. In *Cybersecurity Threats with New Perspectives*. IntechOpen.
- [5] Mario P Brito and Ian GJ Dawson. 2020. Predicting the validity of expert judgments in assessing the impact of risk mitigation through failure prevention and correction. *Risk analysis* 40, 10 (2020), 1928–1943.
- [6] Daniela Soares Cruzes, Martin Gilje Jaatun, Karin Bernsmed, and Inger Anne Tøndel. 2018. Challenges and experiences with applying microsoft threat modeling in agile development projects. In *2018 25th Australasian Software Engineering Conference (ASWEC)*. IEEE, 111–120.
- [7] CyberSeek. 2019. *Cybersecurity Supply/Demand Heat Map*. Retrieved April 21, 2022 from <https://www.cyberseek.org/heatmap.html>
- [8] Chad Dougherty, Kirk Sayre, Robert C Seacord, David Svoboda, and Kazuya Togashi. 2009. *Secure Design Patterns*. Technical Report. Carnegie-Mellon University Pittsburgh, Software Engineering Institute.
- [9] Food and Drug Administration. 2001. Guidance for industry: Statistical approaches to establishing bioequivalence.
- [10] Synopsys Software Integrity Group. 2021. *Building Security In Maturity Model (BSIMM12)*. Retrieved April 20, 2022 from <https://www.bsimm.com>
- [11] Per E Gustafsson. 1998. Gender Differences in risk perception: Theoretical and methodological perspectives. *Risk analysis* 18, 6 (1998), 805–811.
- [12] 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.
- [13] Nasif Imtiaz, Justin Middleton, Joymallya Chakraborty, Neill Robson, Gina Bai, and Emerson Murphy-Hill. 2019. Investigating the effects of gender bias on GitHub. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 700–711.
- [14] Eric Jardine. 2020. The Case against Commercial Antivirus Software: Risk Homeostasis and Information Problems in Cybersecurity. *Risk Analysis* 40, 8 (2020), 1571–1588.
- [15] Johannes G Jaspersen and Gilberto Montibeller. 2015. Probability elicitation under severe time pressure: A rank-based method. *Risk Analysis* 35, 7 (2015), 1317–1335.
- [16] J. Jung, K. Hoefig, D. Domis, A. Jedlitschka, and M. Hiller. 2013. Experimental comparison of two safety analysis methods and its replication. *International Symposium on Empirical Software Engineering and Measurement* (2013), 223–232. <https://doi.org/10.1109/ESEM.2013.59>
- [17] Peter Karpati, Andreas L Opdahl, and Guttorm Sindre. 2011. Experimental comparison of misuse case maps with misuse cases and system architecture diagrams for eliciting security vulnerabilities and mitigations. In *Availability, Reliability and Security (ARES), 2011 Sixth International Conference on*. IEEE, 507–514.
- [18] Peter Karpati, Guttorm Sindre, and Raimundas Matulevicius. 2012. Comparing misuse case and mal-activity diagrams for modelling social engineering attacks. *International Journal of Secure Software Engineering (IJSSSE)* 3, 2 (2012), 54–73.
- [19] Katsiaryna Labunets, Fabio Massacci, and Federica Paci. 2017. On the equivalence between graphical and tabular representations for security risk assessment. In *International Working Conference on Requirements Engineering: Foundation for Software Quality*. Springer, 191–208.
- [20] Katsiaryna Labunets, Fabio Massacci, Federica Paci, et al. 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.
- [21] R Murray Lindsay and Andrew SC Ehrenberg. 1993. The design of replicated studies. *The American Statistician* 47, 3 (1993), 217–228.
- [22] Michael Meyners. 2012. Equivalence tests—A review. *Food quality and preference* 26, 2 (2012), 231–245.
- [23] Andreas L Opdahl and Guttorm Sindre. 2009. Experimental comparison of attack trees and misuse cases for security threat identification. *Information and Software Technology* 51, 5 (2009), 916–932.
- [24] M. Riaz, J. King, J. Slankas, L. Williams, F. Massacci, C. Quesada-López, and M. Jenkins. 2017. Identifying the implied: Findings from three differentiated replications on the use of security requirements templates. *Empirical Software Engineering* 22, 4 (2017), 2127–2178. <https://doi.org/10.1007/s10664-016-9481-1>
- [25] Gema Rodríguez-Pérez, Reza Nadri, and Meiyappan Nagappan. 2021. Perceived diversity in software engineering: a systematic literature review. *Empirical Software Engineering* 26, 5 (2021), 1–38.
- [26] S. Rueda, J.I. Panach, and D. Distant. 2020. Requirements elicitation methods based on interviews in comparison: A family of experiments. *Information and Software Technology* 126 (2020). <https://doi.org/10.1016/j.infsof.2020.106361>
- [27] Per Runeson. 2003. Using students as experiment subjects—an analysis on graduate and freshmen student data. In *Proceedings of the International Conference on Empirical Assessment in Software Engineering*. 95–102.
- [28] Ilaah Salman, Ayse Tosun Misirli, and Natalia Juristo. 2015. Are students representatives of professionals in software engineering experiments?. In *Proceedings of the International Conference on Software Engineering—Volume 1*. IEEE Press, 666–676.
- [29] Adrian Santos, Sira Vegas, Markku Oivo, and Natalia Juristo. 2019. A procedure and guidelines for analyzing groups of software engineering replications. *IEEE Transactions on Software Engineering* 47, 9 (2019), 1742–1763.
- [30] Joanna C. S. Santos, Katy Tarrit, and Mehdi Mirakhorli. 2017. A Catalog of Security Architecture Weaknesses. In *Proceedings of the International Conference on Software Architecture Workshops (ICSAW)*. IEEE Computer Society, 220–223. <https://doi.org/10.1109/ICSAW.2017.25>
- [31] Riccardo Scandariato, Kim Wuyts, and Wouter Joosen. 2015. A descriptive study of Microsoft’s threat modeling technique. *Requirements Engineering* 20, 2 (2015), 163–180.
- [32] DL Schuurmann. 1981. On hypothesis-testing to determine if the mean of a normal-distribution is contained in a known interval. *Biometrics* 37, 3 (1981), 617–617.
- [33] Adam Shostack. 2014. *Threat modeling: Designing for security*. John Wiley & Sons.
- [34] R Stevens, D Votipka, and E.M. Redmiles. 2018. The Battle for New York: A Case Study of Applied Digital Threat Modeling at the Enterprise Level. In *SEC’18: Proceedings of the 27th USENIX Conference on Security Symposium*. USENIX Association, 621–637.

- [35] Jakita O Thomas, Nicole Joseph, Arian Williams, Jamika Burge, et al. 2018. Speaking truth to power: Exploring the intersectional experiences of Black women in computing. In *2018 Research on Equity and Sustained Participation in Engineering, Computing, and Technology (RESPECT)*. IEEE, 1–8.
- [36] Katja Tuma, Gül Calikli, and Riccardo Scandariato. 2018. Threat analysis of software systems: A systematic literature review. *Journal of Systems and Software* 144 (2018), 275–294.
- [37] Katja Tuma, Christian Sandberg, Urban Thorsson, Mathias Widman, Thomas Herpel, and Riccardo Scandariato. 2021. Finding security threats that matter: Two industrial case studies. *Journal of Systems and Software* 179 (2021), 111003.
- [38] Katja Tuma and Riccardo Scandariato. 2018. Two architectural threat analysis techniques compared. In *European Conference on Software Architecture*. Springer, 347–363.
- [39] Katja Tuma, Laurens Sion, Riccardo Scandariato, and Koen Yskout. 2020. Automating the early detection of security design flaws. In *Proceedings of the 23rd ACM/IEEE International Conference on Model Driven Engineering Languages and Systems*. 332–342.
- [40] Dimitri Van Landuyt and Wouter Joosen. 2021. A descriptive study of assumptions in STRIDE security threat modeling. *Software and Systems Modeling* (2021), 1–18.
- [41] Yi Wang and David Redmiles. 2019. Implicit gender biases in professional software development: An empirical study. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Society (ICSE-SEIS)*. IEEE, 1–10.
- [42] George Wright, Fergus Bolger, and Gene Rowe. 2002. An empirical test of the relative validity of expert and lay judgments of risk. *Risk Analysis: An International Journal* 22, 6 (2002), 1107–1122.

Pre-print