

A Comprehensive, Automated Security Analysis of the Uptane Automotive Over-the-Air Update Framework: Technical Report

Robert Lorch
robert-lorch@uiowa.edu
The University of Iowa
Iowa City, Iowa, USA

Cesare Tinelli
cesare-tinelli@uiowa.edu
The University of Iowa
Iowa City, Iowa, USA

Daniel Larraz
daniel-larraz@uiowa.edu
The University of Iowa
Iowa City, Iowa, USA

Omar Chowdhury
omar@cs.stonybrook.edu
Stony Brook University
Stony Brook, New York, USA

ABSTRACT

In this paper, we present our experience of formally verifying the desired security properties of the Uptane over-the-air (OTA) software update framework against a set of applicable threat models. Uptane is gaining traction in the automobile industry and is widely considered the next de-facto standard for OTA automobile software updates. The security of Uptane is of utmost importance because modern automobiles rely on software for their safety-critical functionalities and, especially, require OTA software updates to add new safety features or patch bugs in existing ones. Design flaws in Uptane can either violate the integrity of the updates to be installed or prevent vehicles from installing new updates, both of which can cause severe safety issues. Existing approaches to protocol verification either fail to capture the necessary features of Uptane or suffer from termination issues due to Uptane's complexity. A key component of our approach lies in the *eager* combination of an infinite-state model checker and a cryptographic protocol verifier, where (in contrast to prior *lazy* approaches) we are able to simultaneously eliminate a key manual step in the workflow while enabling reasoning over more fine-grained message structures. In addition, our approach utilizes *two proven soundness- and completeness-preserving state-space-reduction optimizations* for computational tractability, as well as a *meta-level analysis technique* that makes it feasible to reason over Uptane's set of optional protocol features. Our approach is able to discover six new vulnerabilities while rediscovering all five known ones. While there have been previous analyses of Uptane's security properties, they either missed design flaws identified by our approach or suffered from coverage and termination issues. The Uptane standards body has positively acknowledged our findings and has suggested updates to the protocol specification documents to address them.

CCS CONCEPTS

• Security and privacy → Logic and verification.

KEYWORDS

automotive security, model checking, protocol analysis, attacks, vulnerabilities

ACM Reference Format:

Robert Lorch, Daniel Larraz, Cesare Tinelli, and Omar Chowdhury. 2024. A Comprehensive, Automated Security Analysis of the Uptane Automotive Over-the-Air Update Framework: Technical Report. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 24 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Safety-critical functions in modern automotive vehicles such as airbag systems, engine control, stability control, braking, and more are now controlled by software [48]. In fact, the average new car runs on a staggering 100 million lines of code [49]. Therefore, vehicles must regularly perform software updates in order to patch bugs and implement new features. Physically taking vehicles to dealers for software updates can be both ineffective and impractical. It is thus preferable for next-generation automotive vehicles to perform these updates over the air (OTA). OTA updates, however, can be vulnerable to adversaries, as demonstrated by the compromising of many major update repositories including those run by Apache, Debian, Fedora, and GitHub [19, 28, 29, 63]. Further, many major auto manufacturers have suffered security vulnerabilities allowing attackers to tamper with safety-critical functions such as locking/unlocking, engine starting/stopping, and more [16, 18, 64].

Due to these issues, much work has been dedicated to securing OTA updates [2, 4, 10, 11, 14, 21, 30, 34, 36, 41–43, 52, 55, 56, 58]. Deviating from traditional computing devices, software in automotive vehicles runs on electronic control units (ECUs), which have limited computing power, secure remote communication capabilities, and built-in security measures. *An automotive OTA update protocol is thus required to take these limitations into account.*

Uptane [25, 41] is a proposed OTA protocol that is arguably becoming the industry-standard framework for secure automotive OTA updates. It is currently being integrated into Automotive Grade Linux (AGL) [24], which has a huge share of the automotive OS market and is used by manufacturers such as Toyota, Honda, Ford, Nissan, Mercedes, and Audi [24]. In addition to AGL, Uptane

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

© 2024 Association for Computing Machinery.

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

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

is supported by commercial software management platforms (e.g., [1]) with major clients such as Infiniti, Renault, Bosch, and Continental. *Because of this, a large proportion of vehicles on the road will soon be using Uptane.* Design-level protocol issues will become much more problematic after large-scale deployment, as the cost and technical difficulty of fixing protocol design bugs is disproportionately high at later stages of the development cycle. Thankfully, we are currently in a unique window of opportunity to rigorously analyze the protocol for bugs *before* it is widely deployed. *Despite the automotive industry's large reliance on Uptane, its security has not been studied comprehensively in the literature.*

Problem and scope. Uptane claims to protect *safety-critical systems* from adversaries with a wide range of capabilities. Being a framework, Uptane is designed to be compatible with a wide range of automakers and ECU capabilities, meaning that there is considerable freedom for implementers to make design decisions. Since frameworks are specified at a high level of abstraction, design errors can have sweeping consequences. Given this landscape, we answer the following question: *Can we apply formal verification to analyze the Uptane framework with respect to a set of feasible adversaries in order to increase confidence in its promised security properties and discover underlying design issues?* In our formal analysis, we consider adversarial influence both (i) within the vehicle's internal network and (ii) between the vehicle and remote communication partners. *The overall problem can be reduced to a protocol analysis in the symbolic model in which some or all of the above communication channels are controlled by a Dolev-Yao [20] adversary.*

Challenges. A comprehensive, fine-grained, terminating, and automated security analysis of Uptane warrants addressing the following challenges: (1) **long execution traces:** the Uptane protocol induces long execution traces when it is exactly modeled as described in the standard documents, resulting in severe non-termination issues for the analysis; (2) **properties referencing concrete payload values:** analyzing Uptane requires reasoning about fine-grained security properties that reference concrete values in message payloads, which is a problem for both cryptographic protocol verifiers and current abstraction-refinement-based techniques; (3) **large number of problem instances to solve:** when considering Uptane's specific threat models, desired properties, and optional features (i.e., 9 threat models, 8 properties, and 4 optional features), a naive protocol verification approach for *comprehensively* analyzing Uptane requires solving $1152 (= 9 \times 8 \times 2^4)$ protocol verification problem instances; (4) **manual efforts:** although manual efforts needed for formally analyzing a protocol are unavoidable, especially for formalizing the protocol design and collecting desired properties (i.e., inputs to the analysis approach), any protocol analysis approach that requires manual intervention becomes prohibitive when considering 1152 problem instances.

Note that these challenges apply to formal verification in general, not just Uptane. In particular, since automated verification tools are mostly black-box, remedying verification issues (in the form of non-termination, excessive manual effort, etc.) takes expert knowledge, and there is no one-size-fits-all solution. In our analysis of Uptane, we found that none of prior approaches from the literature [31, 32, 37, 38, 41, 47, 50] were sufficient to address the above challenges (see Section 3).

Approach. Our automated analysis approach takes three inputs, namely, (a) a formal protocol model of Uptane, (b) a list of security guarantees formalized as temporal safety properties, and (c) a threat model to consider. It then exhaustively checks whether the design satisfies all the properties with respect to the given threat model. In case of a violation, like any protocol verification approach, it outputs a *counterexample* as an evidence of violation. Our overall approach is aided by the following insights. *Although the insights are demonstrated in terms of analyzing the Uptane protocol, they are equally applicable to the analysis of other protocols.*

□ **A highly automated, staged protocol analysis:** Our staged analysis approach uses the *abstraction-refinement* paradigm in which a cryptographic protocol verifier and a model checker are *eagerly* combined. In contrast to prior *lazy combination approaches* [31, 32], our approach enjoys a higher degree of automation due to the use of an infinite-state model checker and a novel encoding of a *replay attacker*. Both of these allow us to perform a fine-grained reasoning over message payload values while also automatically incorporating the feedback of a cryptographic protocol verifier without *any* manual intervention. A detailed comparison (see Table 2) with prior approaches is discussed in Section 3.

□ **Protocol step and input trace compressions:** Our analysis takes advantage of two optimization techniques, namely, protocol step and input trace compressions. The first of these coalesces different protocol steps into a single step whereas the latter only considers timesteps in which something meaningful happens. Although optimizations like these are sometimes used in formal verification for computational tractability, the differentiating factor in our work is that we have *proven* that our optimizations maintain analysis correctness. That is, the optimized analysis *provably* will not miss any attacks or introduce spurious ones compared to the naive modeling (the proof is included in Appendix D).

□ **Meta-level protocol analysis:** Our approach also enables a meta-level analysis of optional features, greatly reducing the manual efforts to a feasible level. In particular, it allows us to formulate only 72 protocol verification instances (instead of 1152) by considering all possible combinations of the optional features together, which is a $16\times$ reduction. This meta-level analysis is facilitated by a new feature of certain model checkers called *blame assignment* [45, 46].

Prior Analyses of Uptane. In particular, one was performed by [8]—however, our analysis is more comprehensive with respect to the set of threat models and desired security properties. Additionally, the prior work [8] notes that the execution time of their difficult Tamarin lemmas ranged from 20+ hours to non-terminating, while our novel insights into automated protocol verification lead to termination in under 10 minutes in a vast majority of cases. The difference between our approach and [8] is discussed in depth in Section 3 and is outlined in Table 7 in the appendix.

Evaluation and Findings. All of these novel aspects allow us to achieve termination in 69 of 72 attempted protocol verification instances (see Table 4) when comprehensively analyzing the Uptane protocol design. We achieve termination on more instances than other approaches, as discussed later. All our models are publicly available at <https://github.com/lorchrob/UptaneModels>.

Our analysis rediscovered five known attacks, which is an initial validation of the correctness of the approach. In addition, we

discovered six *new* vulnerabilities (see Table 5), confirmed by the standards body, which allow adversaries to prevent vehicles from installing new updates and also cause vehicles to install the wrong updates.

Contributions. In summary, the paper makes the following technical contributions: (i) a *formalization* of the Uptane protocol from a natural-language specification (containing gaps and ambiguities) that goes beyond having only one protocol participant for every role; (ii) six *new security vulnerabilities* in the protocol; (iii) *proofs for the absence of attacks* with respect to our model for 14 combinations of properties and threat models; and, more generally, (iv) a novel *workflow* that enables the faithful modeling of a variety of adversarial scenarios and optional protocol features in a cryptographic setting.

2 PRELIMINARIES

System description. Uptane aims to *protect the ECU software update process*, focusing specifically on the delivery and verification of software updates. Figure 1 gives an example ECU network architecture with two network segments (called *buses*): a Controller Area Network (CAN) bus and a Local Interconnect Network (LIN) bus. Different network segments are connected by *gateway ECUs* that forward traffic from one bus to another.

ECUs updates are usually installed OTA. Since many ECUs lack computing power, only a small number of the ECUs connect to a remote repository over the internet to retrieve new software images for themselves and other ECUs in the vehicle. The new images are distributed to the other ECUs over the buses.

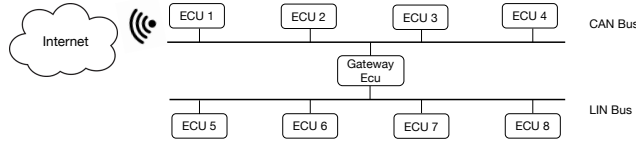


Figure 1: Example ECU Network Architecture

Motivation. Automotive buses were designed with efficiency and convenience in mind, so they *do not include built-in security features* [9]. Therefore, compromising even a single ECU renders the vehicle’s internal network vulnerable to man-in-the-middle attacks. Many techniques to compromise ECUs are *completely remote*. Some common attack vectors are through the OBD-II port, vehicle entertainment systems (e.g. CD drive, USB port), keyless entry systems, and bluetooth and cellular channels [13].

Most effort in automotive bus security has focused on CANs. Efforts in CAN security [5, 23, 40, 51, 53, 54, 59, 60, 66] employ a threat model where the adversary can arbitrarily read, drop, and inject messages in communications between ECUs. Notably, the model allows *ECU impersonation attacks*, where non-safety-critical ECUs can drop CAN packets sent by safety-critical ECUs or inject CAN packets that impersonate safety-critical ECUs [17, 23, 65]. Other studies in automotive bus security are surveyed in [22]— for example, LIN bus message injection attacks involving the error handling mechanism [62]. There are recent news stories of thieves using CAN injection attacks to steal (relatively new) cars [16, 64].

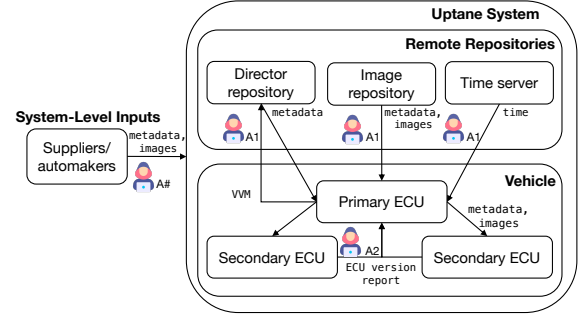


Figure 2: Simplified system architecture of Uptane.

Several studies show realistic man-in-the-middle (MiTM) attacks between an ECU and a communication partner outside the vehicle. For example, [57] demonstrates a remote attacker’s ability to tamper with the tire pressure monitoring system by spoofing ECU sensor data, while [39] involves a remote attacker executing a DoS attack on EV charging through a MiTM attack on the charging protocol.

2.1 Uptane Background

In the Uptane [25, 41] ecosystem (see Figure 2), each vehicle has a *primary ECU* that communicates with two remote repositories, the *director* and *image* repositories. The director repository provides metadata about the software images the ECUs should install. An additional set of metadata, and the images themselves, are retrieved from the image repository. ECUs cross-reference metadata from both repositories, enabling the detection of an attack when one of the repositories is compromised. These checks on the metadata are called *metadata verification*, and they are performed by all ECUs in the vehicle. There are four types of metadata: *root*, *timestamp*, *snapshot*, and *targets*. Some pieces of data are present in *every* metadata file, including a metadata version number and expiration timestamp. For specifics about other information contained in metadata files, consult Appendix B. The primary ECU downloads and verifies images and metadata, which it then distributes to the other ECUs in the vehicle, the *secondary ECUs*. If an ECU is compromised, other ECUs may still resist downloading malicious software, as each ECU performs its own verification. Thus, verification is performed twice on metadata and images for secondary ECUs— once by the primary ECU, and once by each of the secondary ECUs.

Uptane Update Process. An ECU’s software update process can roughly be divided into three stages: (i) update discovery, (ii) metadata verification, and (iii) image verification and acquisition. A high-level Uptane description is given in Figure 8 in the appendix.

Update discovery: In this phase, new images are selected for the ECUs to install. This selection is performed by the director repository, which retrieves information about the vehicle’s currently installed information (aka vehicle version manifest or VVM, for short) and performs dependency resolution to select new updates.

Metadata verification: In this stage, ECUs process metadata from both remote repositories, applying verification procedures to ensure that the metadata is not tampered with.

Image verification and acquisition: In this stage, ECUs apply a verification procedure to the new images retrieved from the image

repository. If this verification is successful, then the ECU installs the new image. If any of the steps above fails, the offending file (metadata, image, or VVM) is discarded and the cycle is restarted with a fresh update discovery. For concrete steps regarding the above stages, consult Appendix B.2.

Metadata verification is flexible for secondary ECUs. Secondary ECUs with less computing power may choose to verify, at a minimum, only targets metadata files from the director repository (*partial verification*), rather than verifying the full set of metadata. In our analysis, we model the entire Uptane system but focus on the security properties of secondary ECUs performing minimal partial verification. Uptane claims to offer *flexibility* without compromising *security*, and we put this claim to the test.

We omitted from the above discussion another type of verification: *verification of the latest time*. Since ECUs often do not have accurate internal clocks, the Uptane system allows for a remote attestation of the latest time, which ECUs must verify before performing metadata verification. However, the specifics of the ECUs' access to time are not discussed in the Uptane specification; instead, the source of time is assumed to be secure. In our analysis, we consider the possibility of a man-in-the-middle attack between the ECU and the (remote) secure source of time.

3 PROBLEM DEFINITION AND CHALLENGES

The underlying problem we solve requires formally analyzing the newest version of the Uptane protocol (version 2.1) with respect to a set of properties and a set of threat models. We discuss them both in this section, together with the technical challenges they pose.

Desired Functional Properties. We consider eight integrity properties. We do not model confidentiality or availability properties as such properties are not promised by Uptane. The first two properties relate to *freshness*, stating that ECUs and remote repositories will not accept old versions of files when newer versions are available.

- P1. ECUs only verify the latest available targets metadata (checks for *freeze attacks*).
- P2. The director repository only verifies the latest available vehicle version manifest (VVM, containing information about each ECU's currently installed image) that was sent by the vehicle (checks for *VVM replay attacks*).

If these properties are violated, the ECUs will continue to run old software images that potentially contain known bugs or security vulnerabilities which could be exploited by an adversary in another attack. Possible implications of the attacks under analysis are outlined in Table 8 in the appendix.

The remaining properties are about ECUs installing the correct software images. Violating them could result in ECU installing (i) older images with known vulnerabilities, (ii) adversary-authored images, or (iii) images with compatibility issues. They can lead to a degradation in vehicle functionality or, in the worst case, adversary control over the vehicle.

- P3. ECUs always verify images in nondecreasing version order (checks for *rollback attacks*).
- P4. ECUs never verify adversary-authored software (checks for *arbitrary software attacks*).

- P5. Remote repositories never verify VVMs containing version reports that were authored by an adversary (checks for *attacker-authored VVM attacks*).
- P6. ECUs never verify metadata instructing installation of an incompatible set of images (checks for *mix-and-match attacks*).
- P7. ECUs always have a compatible set of installed images (checks for *mixed-bundles attacks*).
- P8. ECUs never verify images incompatible with their own hardware (checks for *incompatible image attacks*).

Threat models. We consider an adversary who can read, modify, inject, and drop protocol packets within the vehicle's internal network *and* between the vehicle and remote communication partners, all the while conforming to cryptographic assumptions. This problem reduces to protocol analysis in the symbolic model where selected components (including components within the vehicle) are assumed to be under the influence of a Dolev-Yao[20] adversary.

We consider the following threat models with various adversarial capabilities. Each threat model contains a set of man-in-the-middle attackers (see Figure 2) and access to specified compromised keys.

References to justifications for the technical feasibility of each of these adversarial capabilities are listed in Table 1.

B. The benign case (no adversary).

- A1. The adversary can perform MiTM attacks on network connections *outside* the vehicle (A1 in Figure 2), but not in the vehicle's internal network.
- A2. The adversary can perform MiTM attacks *outside and inside* the vehicle (A1 and A2 in Figure 2).
- A3. As in A2 but the adversary has also access to compromised director repository metadata signing keys.
- A4. As in A3 plus adversary access to compromised image repository timestamp and snapshot metadata signing keys.
- A5. As in A4 plus access to compromised image repository targets metadata signing keys.
- A6. As in A5 plus access to compromised image repository root metadata signing keys.
- A*. As in A1 plus access to the compromised VVM signing key used by the primary ECU.
- A#. The adversary has compromised the supply chain (A# in Figure 2).

Need for Different Threat Models

While the Dolev-Yao [20] model is generic enough to capture the behavior of all of the above threat models, we still must specify which network connections are susceptible to a Dolev-Yao attacker, as well as which cryptographic keys the Dolev-Yao attacker is assumed to have access to.

3.1 Technical challenges

We now illustrate technical challenges by discussing the shortcomings of five existing protocol verification approaches— (1) *manual analysis*, (2) *cryptographic protocol verification in the computational model*, (3) *cryptographic protocol verification in the symbolic model*, (4) *model checking*, and (5) *the abstraction-refinement paradigm*, which lazily combines a symbolic model checker and a symbolic cryptographic protocol verifier.

Adversarial Capability	Feasibility
A1 message reading	[57], [39]
A1 message injection	[57], [39]
A1 message dropping	[57], [39]
A2 message reading	[5, 17, 40, 51, 53, 60, 62, 66]
A2 message injection	[5, 17, 23, 40, 51, 53, 60, 62, 64, 66]
A2 message dropping	[40, 53, 60, 62]
A3-A6, A* key compromise	discussed in Section 9
A# supply chain attack	[15, 35, 61]

Table 1: Feasibility of Adversarial Capabilities. A1 operates outside the vehicle, while A2 operates also inside (Figure 2).

(1) **Manual Analyses.** A naive approach is to manually analyze the protocol to verify its security properties. However, our results (discussed in Section 7) provide strong evidence that manual analysis can miss violations easily identified by automated reasoning.

(2) **Computational Cryptographic Protocol Verifiers.** Protocol verifiers in the computational model (e.g., [3, 6]) represent protocol messages as bitstrings, and cryptographic operators as functions from bitstrings to bitstrings. However, this precise, low-level modeling imposes a cost—high manual effort. For Uptane, this difficulty is exacerbated by the presence of optional features, which exponentially increase the required modeling effort. *In this paper, we instead focus on techniques with high automation.*

(3) **Symbolic Cryptographic Protocol Verifiers.** Protocol verifiers in the symbolic model (e.g., [7, 50]), in short CPVs, assume cryptographic building blocks to be secure (i.e., the *perfect cryptography assumption*). CPVs only reason about the composition of cryptographic primitives, making the protocol verification task more amenable to automation. We tried this approach first, with the Tamarin CPV [50], but we faced two major challenges. The first, which we will call C_1 , is that the analysis failed to terminate. One reason for this is that the search space is large, as Uptane’s update process includes a minimum of 10 steps in total (one for VVM verification, one for each metadata file for each repository, and one for image verification), each relying on ECU or repository state. Specifically, the CPV analysis times out when attempting to prove properties that require reasoning over multiple verification cycles (causing *long execution traces*). The second reason is that Uptane requires fine-grained modeling of infinite-domain data types, arithmetic operations, and temporal operators (e.g., not just standard correspondence and injective-correspondence proofs over symbolic messages). While Tamarin supports these features, the analysis was not scalable enough for Uptane. To address the termination issue, we employed custom heuristics (a Tamarin oracle) to help the CPV’s proof search, but it was not adequate for termination.

The CPV termination issue C_1 is also observed by [8], which is a recent formal analysis of Uptane using Tamarin. In their analysis, powerful compute servers (400GB RAM) were needed, and some proofs still took 20+ hours to complete or failed to terminate. Compared to this approach, ours gives a better performance, even across a wider set of threat models and desired security guarantees

in a commodity laptop. For example, [8] considers a hierarchy of five threat models, where the more powerful adversaries can compromise more components of the model (key compromise is not independently considered). Compared to this, our threat hierarchy of nine threat models is more fine-grained, as it explores all possible key compromises. While it may seem like our threat models do not analyze component compromise, it is equivalently captured by the compromise of all the component’s keys (e.g., threat model A* effectively models a primary ECU compromise). Additionally, [8] considers seven security properties, three of which are from the standard (relating to confidentiality, denial-of-service, and arbitrary software attacks), and four of which they independently formulated (one dealing with privacy, and two dealing with agreement/synchronization between model components). However, two of the properties from the standard (confidentiality and denial-of-service) are trivially broken, and one of the independently formulated properties is very close to protection from freeze attacks (from the standard). In this paper, we consider eight security properties, *not including the two trivial ones*, including five from the standard and three independently formulated. The three that we formulated relate to VVM replay attacks and image installation, and thus have a much higher potential impact if broken (because they relate to update *integrity*). Table 7 in the appendix summarizes the comparison of the two approaches.

Impact of Non-termination

If the analysis does not terminate for a given Uptane problem instance, given by a property, threat model, and a set of optional features, we cannot conclude anything about the property under the threat model for the Uptane protocol design, leaving a potentially undiscovered vulnerability.

Another challenge (C_2) is that Tamarin is unable to efficiently perform a meta-level search over the Uptane framework’s set of optional features. That is, our initial attempt only covered one specific instantiation of the framework. However, we seek to analyze the entire set of all possible (1172) instantiations.

(4) **Model checking.** Another tool-based approach is to use a general-purpose model checker for protocol analysis. General-purpose model checkers excel in fine-grained modeling of infinite-domain data types, arithmetic operations, and temporal operators compared to CPVs. We tried this approach this using Kind 2 [12]. However, Kind 2 (as expected) produces cryptographically infeasible attack traces, as it is not designed to reason about Dolev-Yao adversaries. In addition, Kind 2 had non-termination issues due to long execution traces (C_1).

(5) **Abstraction-Refinement.** A fifth approach is the abstraction-refinement paradigm introduced by Hussain et al. [31, 32], which involves the lazy combination of a cryptographic protocol verifier and a general-purpose model checker. The purpose of the approach is to achieve efficient fine-grained reasoning from a model checker, while reasoning about cryptographic aspects with a CPV. In short, the approach constitutes using a model checker to find an attack trace and then invoking a CPV to determine if the counterexample trace is cryptographically feasible. If it is, the trace returned by the model checker describes a potential attack. If it is not, the

falsified property is updated to block the infeasible counterexample. Unfortunately, even this approach is fundamentally unable to perform a rich, detailed analysis. Its first challenge, denoted by C_3 , is that it relies on predicate abstraction over message payload data. However, in the Uptane use case, an analysis of concrete payload values is required because the desired system-level properties involve arithmetic constraints over packet values (e.g., counters, version numbers, and timestamps). For example, one system-level safety property is protection from rollbacks, formulated as `new_ecu_img.version ≥ curr_ecu_img.version`. Here, `new_ecu_img.version` is part of an incoming message payload and *must* be modeled as an integer to accurately capture the property. This approach also suffers from C_2 , as it does not have support for a meta-level search over optional features in an automated, scalable way. Additionally, the approach does not address C_1 , as the execution traces are still infeasibly long and cause non-termination issues. Finally, this approach warrants an undesirably high number of manual interventions in the workflow C_4 .

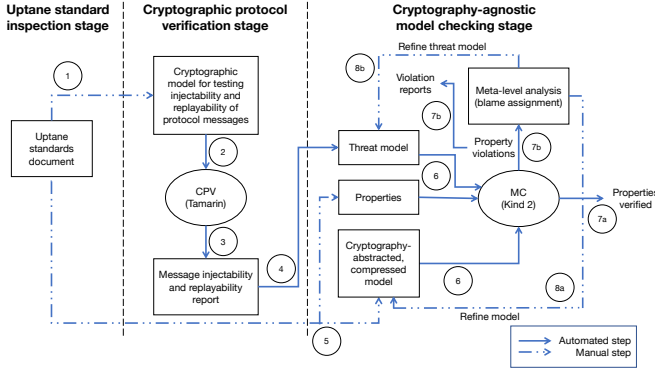


Figure 3: Our Protocol Analysis Workflow (Eager Approach).

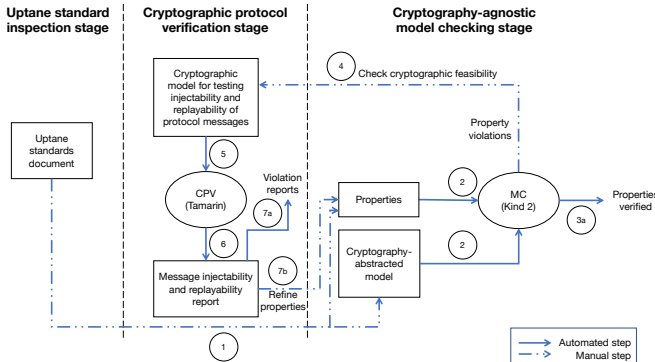


Figure 4: Protocol Analysis Workflow, Lazy Approach [31, 32].

4 APPROACH OVERVIEW

From the above approaches, we consider the abstraction-refinement paradigm to be the most promising as it can reason about both cryptographic and rich protocol features (e.g., statefulness, constraints over payloads). However, we still face four core challenges: long execution traces cause non-termination (C_1); there is no mechanism for meta-level analysis over optional protocol features (C_2); message payload data is modeled with predication abstraction, disallowing analysis of concrete payload values (C_3); and there must be manual intervention for every counterexample trace generated by the model checker (C_4). We will provide the high-level description of our approach and will explain how it is a fundamental improvement over prior approaches [31, 32] using a running example.

4.1 High-level Approach Overview

The high-level approach is outlined in Figure 3. In step ①, we construct a CPV model by consulting the Uptane standards document. In the CPV model, we abstract away control flow and statefulness, only modeling the cryptographic aspects of the protocol. In step ②, we consult CPV to determine the adversary’s cryptographic capabilities, generating a report (step ③). For example, the adversary may not be able to modify messages *arbitrarily*, as this could result in breaking cryptographic assumptions (e.g., an adversary forging a digital signature without knowledge of the corresponding secret signing key). In step ⑤, we create system S and model the desired properties by consulting the Uptane standards document. We model an adversary who modifies messages sent between model components (protocol participants). Based on the CPV report from step ③, we instrument each of S ’s network connections in step ④ by making each connection vulnerable to either (i) no attacks, (ii) only *replay* attacks (i.e., by adversaries that can only replay past messages but not inject new messages), or (iii) arbitrary injections attacks (from unrestricted adversaries). Specifically, if CPV’s corresponding weak authentication lemma is disproven, we consider a standard attacker, where the injection of arbitrary messages is possible. If CPV’s corresponding weak authentication lemma is proven but the injective authentication lemma is disproven, we consider a replay attacker, as replays are feasible but new message injections are not. If CPV’s injective authentication lemma is proven, then we mark the connection as not vulnerable to any attack at all, as even replays are infeasible. The relationship between the CPV lemma results and adversarial capabilities is outlined in Table 3. In step ⑥, we execute the model checker with the corresponding threat model and property. Either the property is verified (step ⑦a), or there is a counterexample trace (step ⑦b).

Contrast Between Lazy and Eager Approaches

The lazy approach (see Figure 4) requires crossing back into the CPV verification stage through arrow 4 (a *manual* backtracking step), which is completely eliminated in the eager approach (see Figure 3).

4.2 Our Approach with Working Example.

In this example, we view our approach as an idealized version of the prior abstraction-refinement paradigm [31, 32]. We will walk through a working example of the approach, being sure to acknowledge whenever one of the challenges ($C_1 - C_4$) arises. In this section, instead of immediately explaining how each challenge is addressed, we will simply assume that there exists some way to overcome it. Then, in the next section, we will outline our novel insights and how they actually enable us to overcome each challenge.

We will walk through a simplified version of the analysis of desired functional requirements P3 and P6, using the idealized workflow. In English, the properties are “ECUs always verify images in nondecreasing version order” and “ECUs never verify metadata that instructs them to install an incompatible set of images,” respectively. Both analyses are primarily concerned with the secondary ECUs’ verification of targets metadata, which is graphically demonstrated in Figure 5. Specifically, Figure 5 demonstrates the high-level steps taken by the various components, numerically indexed to indicate the order of execution. Additionally, the pseudocode describes some of the finer details of each component’s local actions. At a high level, targets metadata is first generated and sent from both repositories to the primary ECU. The ECU performs verification of both metadata files individually, and also performs a cross-reference of the two files to hedge against either repository (or the connection to either repository) being compromised. Finally, the primary ECU sends targets metadata to the secondary ECU, which performs its own verification of metadata. However, in this case, the secondary ECU only processes the director repository’s metadata, so cross-referencing is not possible.

We start with an analysis of P6. In this case, we are assuming threat model A3 (internal adversary with compromised director targets metadata keys).

Create CPV Model. We construct a CPV model by consulting the Uptane standards document.

Execute CPV. For P6, the most relevant lemmas are those corresponding to (i) sending targets metadata from the director repository to the primary ECU (step ①), and (ii) sending targets metadata from the primary ECU to the secondary ECU (step ④). In both cases, the strong authentication lemma fails, meaning that targets metadata replays are cryptographically feasible. However, in case (ii), the weak authentication lemma *also* fails with a counterexample trace where the adversary injects an arbitrary targets metadata file, and it is still verified by the secondary ECU.

Create Model Checker (MC) Model. We create system S by consulting the Uptane standards document. *We assume, in an ideal world, that the model is constructed such that counterexample traces are short and that message payload data is modeled concretely (without predicate abstraction), addressing C_1 and C_3 .*

Specify MC Threat Model. In our example, the results from the CPV model mandate that a *replay attacker* should be placed on the connection from the director to the primary ECU (step ①), and a *standard attacker* should be placed on the connection from the primary ECU to the secondary ECU (step ④).

Execute MC. We execute the model checker with the corresponding threat model. *We assume, in an ideal world, that there is a meta*

Characteristic	Our approach	Hussain <i>et al.</i> [31, 32]	Tamarin [50]
MC + CPV combination strategy	Eager	Lazy	NA
Type of instrumentation	Model instrumentation	Property instrumentation	NA
Concrete message payload?	Yes	No	Yes*
Meta-level analysis (blame assignment)?	Yes	No	No
Compression with meta-theorem?	Yes	No	Yes*

* With scalability issues

Table 2: Comparison of various approaches

analysis which automatically determines which optional protocol features (if any) contributed to the property violation, addressing C_2 .

For P6, MC produces a counterexample (with no optional features disabled) where an adversary injects a targets metadata file on the connection between the primary ECU and the secondary ECU (step ④). This metadata file contains instructions for the ECUs to install an incompatible set of images, representing a violation of P6. *We assume, in an ideal world, that for analyzing P3, we do not have to repeat steps “execute CPV” and “specify MC threat model.” This addresses challenge C_4 .*

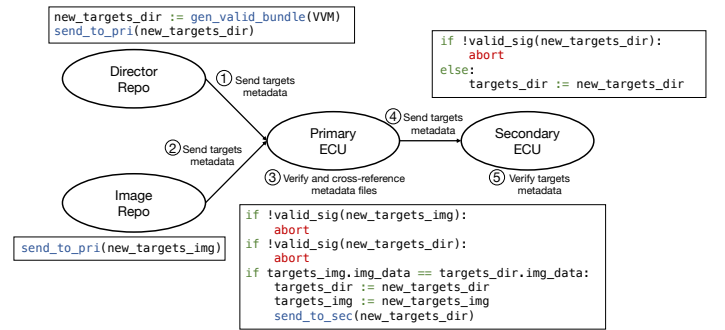


Figure 5: Simplified Targets Metadata Verification Overview.

4.3 Working Example (Prior Work)

We now analyze P3 and P6 again, this time with the approach in [31, 32], to further illustrate the extra manual steps it requires. We start with P6, assuming threat model A3 (internal adversary with compromised director targets metadata keys).

Execute MC. We start with executing the cryptography-agnostic MC model. In this case, MC generates a counterexample trace where a targets metadata file is injected on the connection between the director repository and the primary ECU (step ①).

Execute CPV. We must check the cryptographic feasibility of the attack. In this case, the weak authentication lemma for targets

Relationship between attacker capabilities and Auth. Lemmas	Auth. Lemmas Proven (✓)/ Falsified(×) by CPV	
	Weak Auth. Lemma	Strong Auth. Lemma
Message injection, replay both possible	×	×
Message replay possible, not injection	✓	×
Neither possible	✓	✓

Table 3: Relationship between Authentication Lemmas and Attacker Capabilities (i.e., message injection, message replay)

metadata being sent from the director to the primary is proven (step ①), so the trace is a spurious counterexample.

Property Instrumentation. To address the spurious counterexample, P_6 is instrumented to state “If there is no arbitrary targets metadata injection on the connection from the director to the primary, then ECUs never verify metadata that instructs them to install an incompatible set of images.”

Execute MC. We invoke MC again with the updated property, generating a trace where the targets metadata file is injected on the connection from the primary ECU to the secondary ECU (step ④).

Execute CPV. Again, we check for cryptographic feasibility. Here, the weak authentication lemma for targets metadata sent from the primary ECU to the secondary ECU fails, so the attack is cryptographically feasible. **To analyze P_3 in the lazy approach, the above five steps all have to be repeated.** This is, in part, because the updates to our system model that block cryptographic traces were *property-specific*, i.e., our analysis for P_6 does not carry over to P_3 . Further, this is not the only problem. As discussed in the previous subsection, the abstraction-refinement paradigm still lacks mechanisms for dealing with long execution traces (C_1), performing meta-level analysis over optional protocol features (C_2), and modeling message payload data without predicate abstraction (C_3).

5 TECHNICAL INSIGHTS

In this section, we present an overview of our key technical insights and how they enable the idealized workflow described in the previous section. We outline the differences between our approach and Hussain *et al.* [31, 32] in Table 2.

Insight #1: Compression. To concretely address C_1 , we introduce two compression techniques, *input compression* and *model compression*, which dramatically reduce the length of execution traces. The core ideas are that (1) certain types of system-wide stuttering can be eliminated from modeling, and (2) multiple Uptane protocol steps can be coalesced into a single step while provably preserving semantic correctness with respect to the properties under analysis. In Hussain *et al.*’s approach, the length of attack traces makes analysis intractable. However, with compression, attack traces are limited to a quarter of the length as compared to without compression, leading to termination in more cases (see Table 4). A detailed discussion of compression is in Section 6.

Insight #2: Blame Assignment. To concretely address C_2 , we model optional security-enhancing features symbolically as Boolean parameters in our model and use the *blame assignment* feature of the Kind 2 model checker [45, 46] for fine-grained information about the features that are insufficient for protecting against a certain attack. The blame assignment functionality on a model attempts to find a minimal truth assignment of some given Boolean parameters that is sufficient to trigger a violation of the property. In the truth assignment identified by blame assignment, having value *true* for the Boolean parameter corresponding to an optional feature suggest that the protocol should disable that feature. As a concrete example, consider two optional features in Uptane: (1) Repositories should increment version numbers of metadata files when they are updated, and (2) targets metadata should include image version numbers. With two optional features, there are four concrete protocol instantiations (an example instantiation is to require that version numbers are incremented, but allow targets metadata files to forgo version numbers). In Hussain *et al.*’s approach, the entire workflow has to be replicated for each optional feature (in general, 2^n times for n optional features) to determine which optional features must be disabled for which attacks.

Impact of Blame Assignment

In our approach, a single run of the workflow is comprehensive, as the model checker’s blame assignment feature automatically searches for attack traces that minimize the number of disabled optional features. In our analysis of Uptane, we consider four optional features, reducing manual analysis efforts by a factor of 16.

Insight #3: Replay Attacker. The most common approach of incorporating an abstract version of Dolev-Yao-style adversary is to place a component between every public communication channel that can non-deterministically modify the data being sent through the channel. Without further behavioral restrictions, this adversary actually violates the requirement of the original Dolev-Yao model as it may not conform to cryptographic assumptions (e.g., decrypting a ciphertext without possessing the decryption key). In particular, when the CPV analysis suggests that the attacker can only replay packets but cannot forge new packets, placing a *replay* attacker is more accurate. We designed a non-trivial way to explicitly capture the behavior of a bounded or unbounded replay attacker (see **Insight #6** and Appendix E for more details). The use of the replay attacker helps us address both C_4 and C_3 .

Insight #4: Eager combination with Replay Attacker. To concretely address C_4 , we optimize the workflow from [31, 32] to minimize the number of manual steps. One major difference with those previous works regards the combination of the MC and CPV analyses. The previous works combine them *lazily* as they account for cryptographic influence through *property instrumentation* of the form $\alpha \rightarrow \beta$, where β is the property under consideration and the antecedent α rules out cryptographically infeasible traces. In contrast, we combine the two analyses *eagerly*, accounting for cryptographic influence through *model instrumentation* with the replay attacker.

We demonstrate how we differ with an example. Suppose the existence of four functions, f, g, adv , and $main : \mathbb{Z} \rightarrow \mathbb{Z}$, where we want to prove some properties about the output of $main$:

$$\begin{aligned} f(x) &= |2x| & adv(x) &= ??? \\ g(x) &= 3x & main(x) &= g(adv(f(x))) \end{aligned}$$

We aim to analyze the output of $main$, which composes g and f through adv , where adv models an adversarially-controlled function. In this setting, adv is an *uninterpreted function*, that is, a function whose signature is known but its semantics is not. This simulates the modeling of an adversary whose capabilities are only knowable by consulting an oracle (representing a CPV). Assume in this case that it is only possible for adv to return even integers.

Suppose that we want to verify two properties φ_1 and φ_2 using a suitable automated reasoner (for instance, an SMT solver), where φ_1 states that the output of $main$ is always even and φ_2 states that the output of $main$ is always positive. With no adversarial influence, *i.e.*, when adv is the identity function, both properties trivially hold. In the lazy approach, the behavior of adv is initially unconstrained. So, in the first SMT query for φ_1 , a spurious counterexample is raised where adv returns -1 (or in general, any odd negative number). Then, the counterexample is manually inspected to find that the property failure was due to the injection of an arbitrary value. Only then does the lazy approach invoke the oracle and determine that adv 's behavior was infeasible. Then, φ_1 is instrumented to rule out the attack and rewritten in the form of $\alpha \rightarrow \varphi_1$, where the antecedent α specifies that adv must return an even integer. An inefficiency of this approach is that when querying φ_2 , the same initial counterexample (where adv returns -1) is raised, but the manual work of counterexample inspection and property instrumentation has to be duplicated for φ_2 .

In this paper, we optimize the approach by invoking the CPV ahead of time and specifying adv 's behavior in the model (*not* the property) by giving adv a concrete definition. This minimizes manual trace inspections and completely avoids the issue of having to block similarly infeasible counterexamples multiple times. In other words, the lazy approach (from Hussain *et al.*) involves twice as much manual work in this example compared to our approach—the lazy approach requires *two* CPV invocations and property instrumentations, while the eager approach only requires a single CPV invocation and model instrumentation. If the CPV results denote that protocol components are vulnerable to message replays, then model instrumentation is performed using the replay attacker, discussed previously. (If protocol components are vulnerable to injections of arbitrary messages, we use the standard technique of placing an adversarial component that can arbitrarily, nondeterministically modify messages sent between components, as discussed in **Insight #3**.)

Significance of the Eager Approach

If there are n counterexamples targeting the same cryptographic vulnerability, the eager approach takes n times less manual effort than prior work [31, 32].

Insight #5: Context Independence. We optimize the approach again to address C_4 . We note that in general, we cannot rely on

CPV's output without manually inspecting counterexample traces. For example, if CPV says that a replay attack is possible, then manual analysis of counterexample traces is needed to determine if the attack is *generalizable*, or if it can only take place under specific circumstances. If it is the latter, then the adversary's behavior cannot be accurately modeled with a replay attacker, as additional restrictions are needed on the adversary's behavior. In our analysis of all the CPV lemma counterexamples, we found that all were generalizable, *i.e.*, context independent. This enables us to invoke the CPV only once for each lemma in the workflow, rather than having to revisit the CPV for each discovered attack.

Insight #6: Payloads with Replay Attacker. To concretely address C_3 , we appeal to the *principle of maximum logical revelation*, which states that as many details as possible should be captured in the structure of logical formulae, rather than being abstracted away in the atoms. The replay attacker, described previously, is not restricted to any specific datatype, so it can naturally reason over concrete (*e.g.*, integer) values in message payloads. This was not feasible in [31, 32] because their strategy of accounting for adversarial influence in the system model relies on predicate abstraction.

Consider adding a restriction to our Uptane system model that an arbitrary (previously unseen) ECU version report cannot be injected on the connection from component A to component B . For simplicity, assume the ECU version's payload can be accurately modeled with a pair of natural numbers $\langle a, b \rangle$. In Hussain *et al.*'s approach, the report is abstracted into a Boolean variable v_rep , and a property φ is instrumented to $(Inject(v_rep) \Rightarrow Once(v_rep)) \Rightarrow \varphi$, blocking the injection of v_rep unless it has been sent before. But, their approach does not naturally lift to precisely model v_rep as a pair of natural numbers. Below, suppose v_rep' represents the value of the version report injected by the adversary. Blocking the injection of every pair of natural numbers with property instrumentation is not feasible, as it requires an infinite case split:

$$\begin{aligned} (Inject(v_rep') \wedge v_rep' = \langle 0, 0 \rangle \Rightarrow Once(v_rep = \langle 0, 0 \rangle)) \wedge \\ (Inject(v_rep') \wedge v_rep' = \langle 0, 1 \rangle \Rightarrow Once(v_rep = \langle 0, 1 \rangle)) \wedge \\ \dots \\ \Rightarrow \varphi \end{aligned}$$

*The replay attacker, described previously (**Insight #3**), is thus necessary to model replay attacks generically over infinite domains.*

6 MODEL AND INPUT COMPRESSION

Motivation for Compression. Our model S of the Uptane framework is complex enough that MC is unable to analyze it in reasonable time. The model checker times out for 12 system-level property/threat model combinations. To make the automated analysis of the model feasible, we apply a compression technique to the model that we describe in this section.

We apply two types of compression, *input compression* and *model compression*, graphically illustrated in Figure 6. Intuitively, input compression only considers “meaningful” timesteps with system-level inputs, while model compression coalesces multiple protocol steps to the same step. Both types of compression minimize the

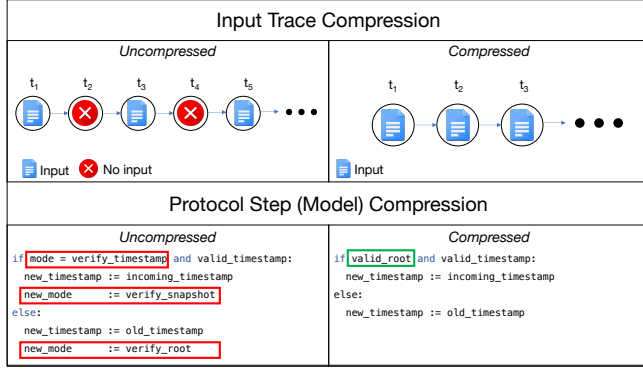


Figure 6: Illustration of Input and Model Compression.

length of verification steps (and hence, attack traces). More details of compression will be discussed later in this section.

Without compression, MC can only verify the absence of attacks in five of the eight properties in the benign case. With compression, all eight are verified, and we also get five additional proofs for the absence of attacks in other threat models. We denote the compressed model as \hat{S} , and we denote S^+ as the model obtained by instrumenting S with Dolev-Yao adversaries at selected communication points. Let \hat{S}^+ denote the instrumented version of \hat{S} .

The results obtained with the compressed model \hat{S} (resp., \hat{S}^+) lift to S (resp., S^+) thanks to a meta-level theorem stating, intuitively, that proofs or violations of properties in \hat{S} correspond to proofs or violations of properties in S . The meta-theorem (Theorem 1) is described graphically in Figure 7. We point out that this correspondence between the compressed and uncompressed system is not a general result; it is restricted to S and the specific integrity properties we consider for Uptane.

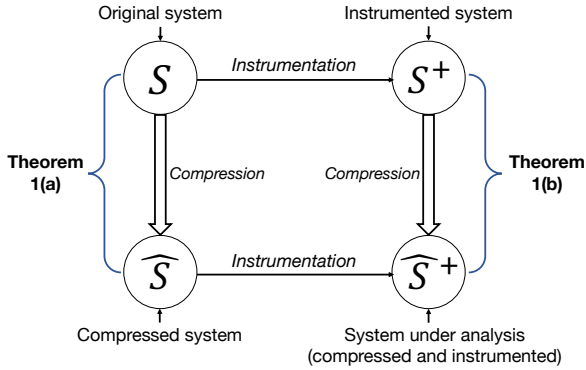


Figure 7: Compression.

Description of Compression. We present more details about S and \hat{S} (the original and compressed systems) and outline the differences between the two. In both S and \hat{S} , the system-level input is comprised of four metadata files (one for each metadata file type) and three software images (one for each ECU). The system-level output is comprised of each ECU's currently verified targets

metadata and currently installed image, as well as the director repository's latest verified VVM. Both the system-level input and system-level output are *event-based*, intuitively meaning that input and output values can be present or not. The absence of a system-level input means that the corresponding image or metadata file is not available for the primary ECU to download.

In S , a full metadata verification cycle spans nine (*execution*) steps where a step corresponds to a transition of the model from a state s to a state s' . Eight steps are each associated with the verification of a single input metadata file, while the ninth step is used to cross-reference metadata from both repositories. Concretely, at every step, the primary ECU is in one of eight (*execution*) modes: (1) waiting for root director metadata, (2) waiting for root timestamp metadata, ..., (8) waiting for image targets metadata, and (9) cross-referencing. At every step, the primary ECU attempts to verify the corresponding metadata file. If it succeeds, then the primary advances to the next mode on the next step. If it fails, the cycle restarts. If there is no incoming metadata file to verify, the ECU considers this a failure and restarts the cycle. If the primary receives an input corresponding to a metadata file that is not currently being verified, it ignores it. Also, images are ignored until the last step of a successful verification cycle. In other words, ECUs do not install new images in the middle of metadata verification. Finally, the primary ECU does not produce output events in the middle of the verification cycle (but only when it starts a fresh verification cycle).

In \hat{S} , the sequential verification steps above in the primary ECU are collapsed to a single step where verification procedures are applied to all input metadata files in parallel. We call this compression *model compression*, to highlight differences in the definition of the system itself in the two models S and \hat{S} . For example, the lower half of Figure 6 demonstrates model compression with timestamp metadata verification. In the uncompressed model, the ECU has to be in the correct mode in order to update its current timestamp metadata. However, in the compressed version, no mode check is required, and the ECU can update a metadata file if it is valid (and all previous metadata files in the verification cycle were valid).

We also perform *input compression*, illustrated in Figure 6, which amounts to considering only execution traces which have a system-level input event at every step. Formally, this is achieved by adding the assumption “always Event” to the set of system-level assumptions where Event is a formula defined to be true in the current state if and only if a system-level input event is present.

Both forms of compressions aim at reducing the length of traces that MC has to reason about to find an *attack trace*, an execution trace that violates one of the properties, or rule out their existence.

The following theorem (which we prove in Appendix D) states, in essence, that applying compression on our model S does not change the satisfiability of the desired properties in either the benign case or with adversarial instrumentation. (In the theorem, \models is the entailment relation in LTL.)

THEOREM 1. *Let A be a formula denoting assumptions on system-level inputs for S . For each desired functional requirement R from Section 2:*

- (a) $S \wedge \text{always } A \models R$ iff $\hat{S} \wedge \text{always } (A \wedge \text{Event}) \models R$
- (b) If $\hat{S}^+ \wedge \text{always } (A \wedge \text{Event}) \not\models R$ iff $S^+ \wedge \text{always } A \not\models R$

7 ANALYSIS FINDINGS

We now discuss the vulnerabilities uncovered by our protocol verification strategy and explain how they could be exploited. We use the Kind 2 [12] model checker to model the non-cryptographic aspects of Uptane and the Tamarin [50] CPV to model the adversary's capabilities with respect to cryptographic assumptions. To reproduce our analysis, please consult Appendix E. Table 4 summarizes the results with respect to property violations, and Table 5 lists the attacks found using our workflow.

Our analysis considers extra attack scenarios A1 and A* which were not present in previous work [37, 38, 47]. For each property, in threat models where no attack was found, we distinguish between MC proving that the property holds and timing out (with \checkmark and $?$). We demonstrate vulnerabilities of three types: (i) showing that previously-known attacks are possible with fewer adversary capabilities; (ii) analyzing attacks that were not previously considered (namely, replayed VVM, attacker-authored VVM, and incompatible image attacks, to be discussed soon); and (iii) demonstrating that the omission of optional features leads to a degradation in security. Our analyses focus on attacks against partial verification secondary ECUs. None of the attack scenarios require the adversary compromising the Primary ECU. We now discuss the new vulnerabilities below. For brevity, we only discuss two findings (**Finding 1** and **Finding OPT1**) in the main text. Details about the rest of the findings are in Appendix F. To further clarify our analysis strategy, the first two findings discussed in the appendix explicitly outline the workflow steps necessary to discover the attacks.

Finding 1: Freeze Attack. The adversary causes a secondary ECU to re-verify the same targets metadata file rather than updating to the latest version.

Threat Model. Internal adversary (A2).

Vulnerability. An adversary (internal or external) can repeatedly replay the same metadata files to the ECU. The ECU checks that the metadata file version number (and for targets metadata, image release counters) are nondecreasing, which will pass. Also, the ECU will check if the metadata file's expiration time is later than the current time. The adversary can block an ECU's access to the current time by dropping the message containing the latest time, either between the external source of time and the primary ECU or between the primary ECU and the secondary ECU. Alternatively, the adversary can inject garbage time messages that will not pass the ECU's verification. The ECU's response to the adversary's tampering with time messages is not well defined in the standard.

In section 5.4.3.1 of the standard, there is no indication to abort the update process if the attestation of the latest time is absent or invalid. This contrasts with sections 5.4.3.2-5.4.3.4, which recommend aborting if metadata or image verification fail. One could interpret this to mean that the verification process should continue *without the ECU updating its clock*, and so the most recent securely attested time would be identical to the last update cycle. This interpretation leads to the freeze attack succeeding.

For more information, one may consult the deployment recommendations [26]. It states in section 3.1.1.1 that the primary ECU should continue without updating its current time if it cannot verify a time message, further supporting the previous paragraph. But,

section 3.1.1.3 says that *all* ECUs should abort, seemingly contradicting section 3.1.1.1 as well as the standards document— more clarification is needed to resolve the situation.

Detection. Based on Tamarin's output, we insert a replay attacker on the connection from the director to the primary and the connection from the primary to the secondary. By assumption, the source of time is secure, so the adversary's only capability is to drop the attestation of the latest time (there are no availability guarantees). This can be modeled by allowing the adversary to replay the time sent along the channel in the previous time step (both from the time server to the primary and from the primary to the secondary).

MC gives a counterexample for the property "When the secondary ECU verifies new targets metadata, the new targets metadata is the latest available." The adversary here drops the message from the primary to the secondary ECU containing the latest time and replays an older targets metadata file. The older targets metadata file is verified by the secondary, even though there is a newer targets metadata file available.

7.1 Analysis of Optional Security Features

One advantage of our approach over [31, 32] is that our approach makes it computationally feasible to reason about complex combinations of optional features. There are some features that are *technically* optional, but their exclusion would lead to immediate and obvious degradation in security. So, rather than disabling *all* optional features, we pick a reasonable subset, that is, we consider a few features whose absence may impact the integrity of the system: (O1) Repositories should increment version numbers of metadata files when they are updated (a repository may update the metadata file but leave the version number unchanged); (O2) targets metadata should include image version numbers (image version numbers are distinct from metadata version numbers); (O3) the director repository should verify digital signatures when checking VVMs for legitimacy; and (O4) the director repository should use information about image dependencies and conflicts when selecting images for vehicles to install. Without using blame assignment, comprehensively analyzing Uptane's security properties with respect to this set of optional features would *require us to execute the entire workflow* $2^4 = 16$ *times*—the blame assignment feature narrows this to just once, greatly reducing manual efforts.

Finding OPT1: Rollback Attack. Disabling just O1 and O2, allows an adversary of weaker capability (A1, compared to A3) to execute a specific kind of rollback attack.

Threat Model. External adversary (A1).

Vulnerability. To verify new targets metadata, the ECU checks the signatures, version number, expiration time, and ECU IDs.

If an adversary replays an older targets metadata file, but the older targets metadata file does not have a lower version number (due to O1), then the version number check will pass. The metadata's signatures and ECU identifiers will still be valid. For the expiration time check, messages containing the latest time can be dropped, or garbage time messages can be injected, as discussed for freeze attacks. This could be performed in the background by the adversary (keeping the ECU's clock at the same value for multiple update cycles) until they spot a vulnerability in one of the ECU's images and are ready to execute the rollback. During image verification,

	Property																							
	P1			P2			P3			P4			P5			P6			P7			P8		
Threat Model	PW	WC	C	PW	WC	C	PW	WC	C	PW	WC	C	PW	WC	C	PW	WC	C	PW	WC	C	PW	WC	C
B	?	✓	✓	NA	✓	✓	?	?	✓	?	✓	✓	NA	✓	✓	?	✓	✓	NA	?	✓	NA	?	✓
A1	NA	×	×	NA	×	×	NA	?	✓	NA	✓	✓	NA	✓	✓	NA	✓	✓	NA	?	✓	NA	?	✓
A2	?	×	×	NA	×	×	?	?	?	?	✓	✓	NA	✓	✓	?	✓	✓	NA	×	×	NA	?	?
A3	?	×	×	NA	×	×	?	×	×	?	×	×	NA	✓	✓	×	×	×	NA	×	×	NA	×	×
A4	×	×	×	NA	×	×	?	×	×	?	×	×	NA	✓	✓	×	×	×	NA	×	×	NA	×	×
A5	×	×	×	NA	×	×	×	×	×	×	×	×	NA	✓	✓	×	×	×	NA	×	×	NA	×	×
A6	×	×	×	NA	×	×	×	×	×	×	×	×	NA	✓	✓	×	×	×	NA	×	×	NA	×	×
A*	NA	×	×	NA	×	×	NA	?	✓	NA	✓	✓	NA	×	×	NA	✓	✓	NA	?	✓	NA	?	?

✓ \mapsto Attack proven impossible with respect to model, × \mapsto Attack found

? \mapsto Attack presumed impossible, but no proof, NA \mapsto Not analyzed

PW \mapsto Prior work[41], WC \mapsto Our approach without compression, C \mapsto Our approach with compression

Green highlight \mapsto Needed compression for termination

Table 4: Partial Verification ECU Security Properties (no optional features omitted)

Attack	Property violated	Our model threat	Prior model threat	Optional features omitted	Contribution
Freeze	P1	$\geq A1$	$\geq A4$		↓
VVM Replay I	P2	$\geq A1$	NA		★
VVM Replay II	P2	$\geq A2$	NA		★
Rollback I	P3	$\geq A3$	$\geq A5$		↓
Rollback II	P3	$\geq A2$	NA	O1, O2	★
Arbitrary software	P4	$\geq A3$	$\geq A5$		↓
Attacker-authored VVM I	P5	$\geq A1$	NA	O3	▽
Attacker-authored VVM II	P5	A*	NA		▽
Mix-and-match	P6	$\geq A3$	$\geq A3$		≡
Mixed-bundles	P7	$\geq A2$	NA		▽
Incompatible image	P8	$\geq A3$	NA		▽

NA \mapsto Not analyzed or not specified, Gray \mapsto no optional features omitted

★ \mapsto New attack, ↓ \mapsto Weaker threat model, ≡ \mapsto Equivalent to prior analysis

▽ \mapsto Mentioned in a previous document, e.g. [26], but there was no analysis of adversarial capabilities necessary for the attack

Table 5: Attacks

the ECU checks for rollbacks by comparing release counters in current and previous targets metadata files. However, this check is not performed if these version numbers are absent (due to the second omission), so the rollback succeeds. An automaker may omit these features if they view two images as essentially equivalent and want to allow an ECU to freely switch between them. This is not intuitively a “rollback,” but the attack demonstrates how important it is to reason about combinations of protocol features.

Detection. Based on output from Tamarin, we insert a replay attacker on the connection from the director repository to the primary and from the primary to the secondary. MC finds a counterexample where the adversary replays a previous targets metadata file with the same version number as the current targets metadata file. (The older targets metadata file contains information about an older image.) The metadata file passes verification. When the primary ECU requests the corresponding image from the image repository, the adversary replays the older image. Image verification passes because the ECU cannot detect that the image is older (because the image’s version number is absent from the targets metadata), and the image matches the ECU’s current metadata.

8 RELATED WORK

There are a few formal analyses of Uptane in the literature, including one [8] which was discussed in depth in Section 3.

Additionally, there is an approach based on Communicating Sequential Processes [37, 38] and an approach based on attack-defense trees [47]. However, these studies have the following limitations: (i) they [37, 38, 47] don’t consider attacks outside of those addressed in [41]; (ii) they [37, 38, 47] don’t consider adversaries operating within the vehicle; (iii) they [37, 38] don’t consider a system with multiple secondary ECUs; (iv) they [37, 38, 47] don’t consider the omission of optional protocol features; (v) they [47] only consider a threat model where the adversary is very powerful (always assuming *all* repository keys are compromised); and (vi) their [37, 38, 47] models did not lead to the discovery of new vulnerabilities. We overcome these limitations and construct a rich model that accounts for multiple different adversarial scenarios and produces counterexamples representing new vulnerabilities.

Previous works with the most similar methods are LTEInspector [31] and 5GReasoner [32], which combine model checkers with

CPVs to analyze 4G and 5G cellular protocols. Our approach builds on these methods in several ways, as discussed in Section 5.

There is also previous work related to OTA updates in general—Uptane is built on top of a general framework called The Update Framework, which is discussed and analyzed in [10, 11, 42, 43, 58]. In addition, many works analyze OTA updates in the context of IoT systems [4, 21, 30, 52], which have similar limitations as automotive systems (e.g., devices having limited computing power). Finally, some works analyze OTA updates in vehicles [2, 14, 34, 36, 55, 56], but consider approaches separate from the Uptane protocol.

9 DISCUSSION

Responsible Disclosure. We reported our findings to the Uptane standards body, and the legitimacy of our results has been positively acknowledged. To address the issues discussed in this paper, we are actively collaborating with the Uptane standards body, including working on updates to the core specification documents [26, 27]. For example, since reporting our findings, both VVM replay attacks (see Appendix F) have been addressed in the newest version of Uptane (version 2.1), and there is an active GitHub issue addressing the freeze attack. Further updates are currently in progress.

Performance. Our Uptane system model consists of around 2K lines of specs and eight functional properties. In table 4, to prove properties and to find counterexamples for disproven properties, MC took 7 minutes and 49 seconds and 2 minutes and 33 seconds on average, respectively. Out of 20 instances where no attack was found, MC was able to prove the absence of an attack in 14 instances. In the other 6 instances, MC timed out. (Performance metrics were gathered using an Apple M2 CPU and 16 GB RAM.)

Feasibility of Threat Models. Threat models A3 through A* rely on various *key compromises*, which are strong assumptions. However, the Uptane standard and deployment recommendations state that director repository keys ought to be kept *online* to automate creation of fresh metadata. In contrast, image repository keys ought to be kept *offline*, as the corresponding metadata is updated less frequently. Note that *all* of our attacks only require compromising *online* director keys, which are much more vulnerable.

Generalizability of Insights. While Insight #1 (compression) is Uptane-specific, all five other insights (comprising the automated verification workflow, including eager combination with replay attackers and blame assignment) are *protocol-agnostic*.

Testbed. To the best of our knowledge, there are no *complete* and *freely available* implementations of Uptane at the time of writing this paper. For example, the reference implementation [44] is no longer mentioned on Uptane website, as it is based on an obsolete version of the standard. Also, the implementation in AGL omits core aspects of the standard, including parts of metadata verification. While experimentation in a testbed would be ideal, we argue that specification-level issues and ambiguities can lead to vulnerable implementations and should be resolved.

Threat to Validity. We put forth a good-faith attempt to model the standard to the best of our ability. Implementations that interpret the standard differently, deviate from the standard, or provide additional security measures not prescribed by the standard may suffer from a different set of vulnerabilities.

10 CONCLUSION

We present a novel workflow for the formal analysis of security protocols and apply it to Uptane, a state-of-the-art OTA update protocol. Our strategy leverages the strengths of model checkers and CPVs, and its application to Uptane reveals five known and six new security vulnerabilities. Since previous manual efforts failed to document these vulnerabilities, we argue that automated reasoning is a uniquely comprehensive and rigorous tool for protocol analysis.

REFERENCES

- [1] Airbiquity. 2024. Airbiquity. <https://www.airbiquity.com/>. Accessed: March 25, 2023.
- [2] N. Asokan, Thomas Nyman, Norrathep Rattanavipanon, Ahmad-Reza Sadeghi, and Gene Tsudik. 2018. ASSURED: Architecture for Secure Software Update of Realistic Embedded Devices. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 11 (2018), 2290–2300. <https://doi.org/10.1109/TCAD.2018.2858422>
- [3] Gilles Barthe, François Dupressoir, Benjamin Grégoire, César Kunz, Benedikt Schmidt, and Pierre-Yves Strub. 2012. EasyCrypt: A tutorial. *International School on Foundations of Security Analysis and Design* (2012), 146–166.
- [4] Jan Bauwens, Peter Ruckebusch, Spiliotis Giannoulis, Ingrid Moerman, and Eli De Poorter. 2020. Over-the-Air Software Updates in the Internet of Things: An Overview of Key Principles. *IEEE Communications Magazine* 58, 2 (2020), 35–41. <https://doi.org/10.1109/MCOM.001.1900125>
- [5] Rohit Bhatia, Vireshwar Kumar, Khaled Serag, Z. Berkay Celik, Mathias Payer, and Dongyan Xu. 2021. Evading Voltage-Based Intrusion Detection on Automotive CAN. In *Network and Distributed System Security Symposium*.
- [6] Bruno Blanchet. 2007. CryptoVerif: Computationally sound mechanized prover for cryptographic protocols. In *Dagstuhl seminar “Formal Protocol Verification Applied*, Vol. 117. 156.
- [7] Bruno Blanchet. 2016. Modeling and Verifying Security Protocols with the Applied Pi Calculus and ProVerif. *Found. Trends Priv. Secur.* 1, 1–2 (Oct. 2016), 1–135. <https://doi.org/10.1561/33000000004>
- [8] Ioana Boureanu. 2023. Formally Verifying the Security and Privacy of an Adopted Standard for Software-Update in Cars: Verifying Uptane 2.0. In *IEEE SMC 2023*.
- [9] Mehmet Bozdağ, Mohammad Samie, and Ian Jennions. 2018. A survey on can bus protocol: Attacks, challenges, and potential solutions. In *2018 International Conference on Computing, Electronics and Communications Engineering (ICCECE)*. IEEE, 201–205.
- [10] Justin Cappos, Justin Samuel, Scott Baker, and John H. Hartman. 2008. A Look in the Mirror: Attacks on Package Managers. In *Proceedings of the 15th ACM Conference on Computer and Communications Security (Alexandria, Virginia, USA) (CCS ’08)*. Association for Computing Machinery, New York, NY, USA, 565–574. <https://doi.org/10.1145/1455770.1455841>
- [11] Justin Cappos, Justin Samuel, Scott Baker, and John H Hartman. 2008. Package management security. *University of Arizona Technical Report* (2008), 08–02.
- [12] Adrien Champion, Alain Mebsout, Christoph Stickel, and Cesare Tinelli. 2016. The Kind 2 Model Checker. In *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 9780)*, Swarat Chaudhuri and Azadeh Farzan (Eds.). Springer, 510–517. https://doi.org/10.1007/978-3-319-41540-6_29
- [13] Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, Stefan Savage, Karl Koscher, Alexei Czeskis, Franziska Roesner, and Tadayoshi Kohno. 2011. Comprehensive Experimental Analyses of Automotive Attack Surfaces. In *20th USENIX Security Symposium (USENIX Security 11)*. USENIX Association, San Francisco, CA. <https://www.usenix.org/conference/usenix-security-11/comprehensive-experimental-analyses-automotive-attack-surfaces>
- [14] Thomas Chowdhury, Eric Lesiuta, Kerianne Rikley, Chung-Wei Lin, Eunsuk Kang, Baekgyu Kim, Shinichi Shiraishi, Mark Lawford, and Alan Wassyng. 2018. Safe and Secure Automotive Over-the-Air Updates. In *Computer Safety, Reliability, and Security*, Barbara Gallina, Amund Skavhaug, and Friedemann Bitsch (Eds.). Springer International Publishing, Cham, 172–187.
- [15] Catalin Cimpanu. 2017. Petya Ransomware Outbreak Originated in Ukraine via Tainted Accounting Software. Available at <https://www.bleepingcomputer.com/news/security/petya-ransomware-outbreak-originated-in-ukraine-via-tainted-accounting-software/>.
- [16] Lindsay Clark. 2023. CAN do attitude: How thieves steal cars using network bus. https://www.theregister.com/2023/04/06/can_injection_attack_car_theft. Accessed: March 25, 2023.
- [17] Gianpiero Costantino and Ilaria Matteucci. 2020. KOFFEE-Kia OFFensive Exploit. *Istituto di Informatica e Telematica, Tech. Rep.* (2020).
- [18] Sam Curry. 2023. Web Hackers vs. The Auto Industry. <https://samcurry.net/web-hackers-vs-the-auto-industry/>. Accessed: March 25, 2023.

- [19] Debian. 2003. Debian Investigation Report after Server Compromises. <https://www.debian.org/News/2003/20031202>. Accessed: March 25, 2023.
- [20] D. Dolev and A. Yao. 1983. On the security of public key protocols. *IEEE Transactions on Information Theory* 29, 2 (1983), 198–208. <https://doi.org/10.1109/TIT.1983.1056650>
- [21] Saad El Jaouhari and Eric Bouvet. 2022. Secure firmware Over-The-Air updates for IoT: Survey, challenges, and discussions. *Internet of Things* 18 (2022), 100508. <https://doi.org/10.1016/j.iot.2022.100508>
- [22] Zeinab El-Rewini, Karthikeyan Sadatsharan, Daisy Flora Selvaraj, Siby Jose Plathottam, and Prakash Ranganathan. 2020. Cybersecurity challenges in vehicular communications. *Vehicular Communications* 23 (2020), 100214. <https://doi.org/10.1016/j.vehcom.2019.100214>
- [23] Ian Foster, Andrew Prudhomme, Karl Koscher, and Stefan Savage. 2015. Fast and Vulnerable: A Story of Telematic Failures. In *9th USENIX Workshop on Offensive Technologies (WOOT 15)*. USENIX Association, Washington, D.C. <https://www.usenix.org/conference/woot15/workshop-program/presentation/foster>
- [24] Linux Foundation. 2024. Automotive Grade Linux. <https://www.automotivelinux.org/>. Accessed: Mar 15, 2023.
- [25] Linux Foundation. 2024. Uptane – Securing Software Updates for Automobiles. <https://uptane.github.io/>. Accessed: Mar 15, 2023.
- [26] Linux Foundation. 2024. Uptane Deployment Best Practices v.2.0.0. https://uptane.github.io/papers/V2.0.0_uptane_deploy.html. Accessed: Mar 15, 2023.
- [27] Linux Foundation. 2024. Uptane Standard for Design and Implementation 2.0.0. <https://uptane.github.io/papers/uptane-standard.2.0.0.html>. Accessed: Mar 15, 2023.
- [28] P. W. Frields. 2008. Infrastructure report, 2008-08-22 UTC 1200. <https://listman.redhat.com/archives/fedora-announce-list/2008-August/msg00012.html>. Accessed: March 25, 2023.
- [29] Inc. GitHub. 2012. Public Key Security Vulnerability and Mitigation. <https://github.blog/2012-03-04-public-key-security-vulnerability-and-mitigation/>. Accessed: March 25, 2023.
- [30] Xinchu He, Sarra Alqahtani, Rose Gamble, and Mauricio Papa. 2019. Securing Over-The-Air IoT Firmware Updates Using Blockchain. In *Proceedings of the International Conference on Omni-Layer Intelligent Systems (Crete, Greece) (COINS '19)*. Association for Computing Machinery, New York, NY, USA, 164–171. <https://doi.org/10.1145/3312614.3312649>
- [31] Syed Rafiul Hussain, Omar Chowdhury, Shagufta Mehnaz, and Elisa Bertino. 2018. LTEInspector: A Systematic Approach for Adversarial Testing of 4G LTE. In *Proceedings 2018 Network and Distributed System Security Symposium*. Internet Society, San Diego, CA. <https://doi.org/10.14722/ndss.2018.23313> tex.ids: hussainLTEInspectorSystematicApproach2018a.
- [32] Syed Rafiul Hussain, Mitziu Echeverria, Imtiaz Karim, Omar Chowdhury, and Elisa Bertino. 2019. 5GReasoner: A Property-Directed Security and Privacy Analysis Framework for 5G Cellular Network Protocol. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (London, United Kingdom) (CCS '19)*. Association for Computing Machinery, New York, NY, USA, 669–684. <https://doi.org/10.1145/3319535.3354263>
- [33] Michael Huth and Mark Ryan. 2004. *Logic in Computer Science: Modelling and reasoning about systems*. Cambridge university press.
- [34] Muhammad Sabir Idrees, Hendrik Schweppe, Yves Roudier, Marko Wolf, Dirk Scheuermann, and Olaf Henniger. 2011. Secure Automotive On-Board Protocols: A Case of Over-the-Air Firmware Updates. In *Communication Technologies for Vehicles*, Thomas Strang, Andreas Festag, Alexey Vinel, Rashid Mehmood, Cristina Rico Garcia, and Matthias Röckl (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 224–238.
- [35] Swati Khandelwal. 2018. CCleaner Attack Timeline—Here's How Hackers Infected 2.3 Million PCs. Available at <https://thehacknews.com/2018/04/ccleaner-malware-attack.html>.
- [36] Muzaffar Khurram, Hemanth Kumar, Adi Chandak, Varun Sarwade, Nitu Arora, and Tony Quach. 2016. Enhancing connected car adoption: Security and over the air update framework. In *2016 IEEE 3rd World Forum on Internet of Things (WF-IoT)*. 194–198. <https://doi.org/10.1109/WF-IoT.2016.7845430>
- [37] Rhys Kirk, Hoang Nga Nguyen, Jeremy Bryans, Siraj Shaikh, David Evans, and David Price. 2021. Formalising UPTANE in CSP for Security Testing. In *2021 IEEE 21st International Conference on Software Quality, Reliability and Security Companion (QRS-C)*. 816–824. <https://doi.org/10.1109/QRS-C55045.2021.00124>
- [38] Rhys Kirk, Hoang Nga Nguyen, Jeremy Bryans, Siraj Ahmed Shaikh, and Charles Wartner. 2023. A formal framework for security testing of automotive over-the-air update systems. *Journal of Logical and Algebraic Methods in Programming* 130 (2023), 100812. <https://doi.org/10.1016/j.jlamp.2022.100812>
- [39] Sebastian Köhler, Richard Baker, Martin Strohmeier, and Ivan Martinovic. 2022. Brokenwire: Wireless disruption of ccs electric vehicle charging. *arXiv preprint arXiv:2202.02104* (2022).
- [40] Karl Koscher, Alexei Czeskis, Franziska Roesner, Shwetank Patel, Tadayoshi Kohno, Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, and Stefan Savage. 2010. Experimental Security Analysis of a Modern Automobile. In *2010 IEEE Symposium on Security and Privacy*. 447–462. <https://doi.org/10.1109/SP.2010.34>
- [41] Trishank Karthik Kuppusamy, Lois Anne DeLong, and Justin Cappos. 2018. Uptane: Security and Customizability of Software Updates for Vehicles. *IEEE Vehicular Technology Magazine* 13, 1 (2018), 66–73. <https://doi.org/10.1109/MVT.2017.2778751>
- [42] Trishank Karthik Kuppusamy, Vladimir Diaz, and Justin Cappos. 2017. Mercury: Bandwidth-Effective Prevention of Rollback Attacks Against Community Repositories. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. USENIX Association, Santa Clara, CA, 673–688. <https://www.usenix.org/conference/atc17/technical-sessions/presentation/kuppusamy>
- [43] Trishank Karthik Kuppusamy, Santiago Torres-Arias, Vladimir Diaz, and Justin Cappos. 2016. Diplomat: Using Delegations to Protect Community Repositories. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*. USENIX Association, Santa Clara, CA, 567–581. <https://www.usenix.org/conference/nsdi16/technical-sessions/presentation/kuppusamy>
- [44] NYU Secure Systems Lab. 2023. Uptane Reference Implementation. <https://github.com/uptane/obsolete-reference-implementation>. Accessed: March 25, 2023.
- [45] Daniel Larraz, Mickaël Laurent, and Cesare Tinelli. 2021. Merit and blame assignment with Kind 2. In *Formal Methods for Industrial Critical Systems: 26th International Conference, FMICS 2021, Paris, France, August 24–26, 2021, Proceedings 26*. Springer, 212–220.
- [46] Daniel Larraz and Cesare Tinelli. 2023. Finding Locally Smallest Cut Sets Using Max-SMT. *ACM SIGAda Ada Letters* 42, 2 (Apr 2023), 32–39. <https://doi.org/10.1145/3591335.3591337>
- [47] Shahid Mahmood, Alexy Fouillade, Hoang Nga Nguyen, and Siraj A. Shaikh. 2020. A Model-Based Security Testing Approach for Automotive Over-The-Air Updates. In *2020 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*. 6–13. <https://doi.org/10.1109/ICSTW50294.2020.00019>
- [48] Farhad Manjoo. 2010. I'm Sorry, Dave, I'm Afraid I Can't Make a U-Turn. <https://slate.com/technology/2010/02/should-we-be-worried-that-our-cars-are-controlled-by-software.html>. Accessed: Mar 15, 2023.
- [49] Anthony Martin. 2020. Vehicle Dynamics International. <https://www.vehicle-dynamics-international.com/features/vehicle-cybersecurity-control-the-code-control-the-road.html>. Accessed: Mar 15, 2023.
- [50] Simon Meier, Benedikt Schmidt, Cas Cremers, and David Basin. 2013. The TAMARIN Prover for the Symbolic Analysis of Security Protocols. In *Computer Aided Verification*, Natasha Sharygina and Helmut Veith (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 696–701.
- [51] Charlie Miller and Chris Valasek. 2015. Remote exploitation of an unaltered passenger vehicle. *Black Hat USA 2015*, S 91 (2015), 1–91.
- [52] Imanol Mugarza, Jose Luis Flores, and Jose Luis Montero. 2020. Security Issues and Software Updates Management in the Industrial Internet of Things (IIoT) Era. *Sensors* 20, 24 (2020). <https://doi.org/10.3390/s20247160>
- [53] Sen Nie, Ling Liu, and Yuefeng Du. 2017. Free-fall: Hacking tesla from wireless to can bus. *Briefing, Black Hat USA 25* (2017), 1–16.
- [54] Sen Nie, Ling Liu, Yuefeng Du, and Wenkai Zhang. 2018. Over-the-air: How we remotely compromised the gateway, BCM, and autopilot ECUs of Tesla cars. *Briefing, Black Hat USA* (2018), 1–19.
- [55] Dennis K. Nilsson, Lei Sun, and Tatsuo Nakajima. 2008. A Framework for Self-Verification of Firmware Updates over the Air in Vehicle ECUs. In *2008 IEEE Globecom Workshops*. 1–5. <https://doi.org/10.1109/GLOCOMW.2008.ECP.56>
- [56] Anam Qureshi, Murk Marvi, Jawwad Ahmed Shamsi, and Adnan Ajiz. 2022. eUF: A framework for detecting over-the-air malicious updates in autonomous vehicles. *Journal of King Saud University - Computer and Information Sciences* 34, 8, Part A (2022), 5456–5467. <https://doi.org/10.1016/j.jksuci.2021.05.005>
- [57] Ishtiaq Rouf, Robert D Miller, Hossen A Mustafa, Travis Taylor, Sangho Oh, Wenyuan Xu, Marco Gruteser, Wade Trappe, and Ivan Seskar. 2010. Security and Privacy Vulnerabilities of In-Car Wireless Networks: A Tire Pressure Monitoring System Case Study. In *USENIX Security Symposium*, Vol. 10.
- [58] Justin Samuel, Nick Mathewson, Justin Cappos, and Roger Dingledine. 2010. Survivable Key Compromise in Software Update Systems. In *Proceedings of the 17th ACM Conference on Computer and Communications Security (Chicago, Illinois, USA) (CCS '10)*. Association for Computing Machinery, New York, NY, USA, 61–72. <https://doi.org/10.1145/1866307.1866315>
- [59] Khaled Serag, Rohit Bhatia, Akram Faqih, Muslum Ozgur Ozmen, Vireshwar Kumar, Z. Berkay Celik, and Dongyan Xu. 2023. ZBCAN: A Zero-Byte CAN Defense System.
- [60] Khaled Serag, Rohit Bhatia, Vireshwar Kumar, Z Berkay Celik, and Dongyan Xu. 2021. Exposing New Vulnerabilities of Error Handling Mechanism in CAN. In *USENIX Security Symposium*. 4241–4258.
- [61] SolarWinds. 2019. SolarWinds Security Advisory. Available at <https://www.solarwinds.com/securityadvisory>.
- [62] Junko Takahashi, Yosuke Aragane, Toshiyuki Miyazawa, Hitoshi Fuji, Hirofumi Yamashita, Keita Hayakawa, Shintaro Ukai, and Hiroshi Hayakawa. 2017. Automotive attacks and countermeasures on lin-bus. *Journal of Information Processing* 25 (2017), 220–228.

- [63] Apache Infrastructure Team. 2009. apache.org incident report for 8/28/2009. <https://blogs.adobe.com/conversations/2012/09/adobe-to-revoke-code-signing-certificate.html>. Accessed: March 25, 2023.
- [64] Ken Tindell. 2023. CAN Injection: keyless car theft. <https://kentindell.github.io/2023/04/03/can-injection/>. Accessed: March 25, 2023.
- [65] Haohuang Wen, Qi Alfred Chen, and Zhiqiang Lin. 2020. Plug-N-Pwned: Comprehensive Vulnerability Analysis of OBD-II Dongles as A New Over-the-Air Attack Surface in Automotive IoT. In *29th USENIX Security Symposium (USENIX Security 20)*. USENIX Association, 949–965. <https://www.usenix.org/conference/usenixsecurity20/presentation/wen>
- [66] Samuel Woo, Hyo Jin Jo, and Dong Hoon Lee. 2014. A practical wireless attack on the connected car and security protocol for in-vehicle CAN. *IEEE Transactions on intelligent transportation systems* 16, 2 (2014), 993–1006.

A ABBREVIATIONS

Table 6 contains the abbreviations we use throughout the paper and their meanings.

Abbreviation	Meaning
MC	model checker
CPV	cryptographic protocol verifier
S	Uptane model
<i>Properties</i>	
P1	checks for freeze
P2	checks for VVM replay
P3	checks for rollback
P4	checks for arbitrary software
P5	checks for attacker-authored VVM
P6	checks for mix-and-match
P7	checks for mixed-bundles
P8	checks for incompatible image
<i>Threat models</i>	
B	benign
A1	outside vehicle MitM
A2	outside + inside vehicle MitM
A3	Full MitM + director keys
A4	A3 + timestamp & snapshot image repo keys
A5	A4 + targets image repo keys
A6	A5 + root image repo keys
A*	A1 + primary ECU key

Table 6: Abbreviations and their Meaning

B UPTANE BACKGROUND

B.1 Uptane Metadata

Uptane uses four types of metadata: *root*, *timestamp*, *snapshot*, and *targets* metadata, each requiring a threshold of digital signatures. Root metadata specifies the public keys associated with each metadata type (root of trust). Targets metadata includes information about images to be installed (image filenames, hashes, etc.). Snapshot metadata contains a filename and version number of each targets metadata file currently on the repository (to guarantee consistency). Timestamp metadata includes the filename and version number of the latest snapshot metadata file, as well as its hash (to guarantee freshness). Each metadata file also has a version number and expiration timestamp. At construction time, the vehicle is manually initialized with a valid set of metadata.

Update discovery: In the *update discovery* phase, new images are selected for the ECUs to install. This selection is performed by the director repository, which retrieves information about the vehicle’s currently installed information (aka vehicle version manifest or VVM, for short) and performs dependency resolution to select new updates.

- 1) The primary ECU queries the director repository for new updates. This query includes information about the vehicle’s currently installed images in a file called the *vehicle version manifest* (VVM), signed with the primary ECU’s secret key.
- 2) The director repository verifies the VVM by checking digital signatures, ECU IDs, and nonces.
- 3) The director repository generates a fresh set of metadata on the new updates to be installed and sends it to the primary ECU.

B.2 Protocol Steps

Metadata verification: In the *metadata verification* stage, ECUs process metadata from both remote repositories, applying verification procedures to ensure that the metadata is not tampered with.

- 4) The primary ECU verifies each metadata file from the director repository. If verification succeeds, the primary ECU overwrites its current set of metadata. We will not discuss the specifics of metadata verification; for more information see [27, 41].
- 5) The primary ECU queries the image repository for fresh metadata and performs the same verification.
- 6) The primary ECU cross-references the targets metadata from the two repositories to check for discrepancies.
- 7) If the verification succeeds, the primary ECU sends the new metadata to all secondary ECUs, which then perform their own metadata verification.

Image verification and acquisition: In the *image verification and acquisition* stage, ECUs apply a verification procedure to the new images retrieved from the image repository. If this verification is successful, then the ECU installs the new image.

- 8) The primary ECU downloads new images from the image repository which correspond to those listed in its latest set of metadata.
- 9) The primary ECU verifies the images by comparing the images’ hashes to those specified in the targets metadata.
- 10) If verification succeeds, the primary installs its new image and forwards the other new images to the secondary ECUs, which also perform verification and installation.
- 11) Each secondary ECU reports back to the primary whether or not the updates were successfully installed in an *ECU version report*. Each report is signed with the corresponding secondary ECU’s secret key and is included in the VVM at the beginning of the next update cycle (update discovery).

If any of the steps above fails, the offending file (metadata, image, or VVM) is discarded and the cycle is restarted with a fresh update discovery.

Attack Type	Our Approach	Approach from[8]
Number of threat models	9	5
Number of security properties	8	7
Meta-level analysis of optional protocol features with blame assignment?	Yes	No
Compression?	Yes	No
Eager combination with replay attacker?	Yes	No

Table 7: Comparison of approach with [8].

B.3 Alice-Bob Description of Uptane

We give a high-level description of the operations performed and the messages sent in a successful Uptane update cycle using Alice-Bob notation (see Figure 8). In the figure, the protocol steps are numerically labeled to illustrate which step of the three stages (outlined in the main text) they correspond to. We include four principles, P (primary ECU), S (secondary ECU), D (director repository), and I (image repository). We denote keys in the form key_{dr} , where (for example) key denotes a secret key, d denotes the Director principal, and r denotes the root role. Local operations/checks are encoded by the *check* keyword. The Alice-Bob notation assumes a “happy path” where all checks succeed. For simplicity of presentation, we assume that the threshold number of signatures required for each metadata file is one and that there is a single secondary ECU. Further, we leave some of the details implicit: (i) that all signatures are verified according to the latest root metadata, and (ii) that each ECU checks that each incoming metadata file has a nondecreasing version number and is not expired.

C PROPERTY FORMULATION

```

P1. always (
    Event^c => targets_meta_secondary
              = director_latest_targets
)

P2. always (
    Event^c => director_latest_manifest
              = primary_latest_manifest
)

P3. always (
    Event^c => image_secondary_version
              >= old_image_secondary_version
)

P4. always (
    Event^c => AuthoredByOem(image_secondary)
)

P5. always (
    Event^c =>
    AuthoredByEcu(director_latest_manifest)
)

```

```

P6. always (
    Event^c => Compatible(targets.image_one,
                          targets.image_two,
                          targets.image_three)
)

P7. always (
    Event^c =>
    (Historically(pri_soft_version
                  = sec_soft_version) =>
    Compatible(pri_image, sec_image, sec2_image)
)

P8. always (
    Event^c =>
    CompatibleHardware(image_sec, hardware_sec)
)

```

D COMPRESSION PROOF

In this section, we discuss and prove Theorem 1 by proving some intermediate results.

Preliminaries. We rely on standard concepts and notation from (First-Order) Linear Temporal Logic (see, e.g., [33]). We will call an LTL formula a *state property* if it contains no temporal operators. We fix a finite set \mathcal{V} of variables of various types, including Boolean and integer and pairs of integers and booleans. We partition \mathcal{V} into three set respectively of *input*, *output* and *internal* variables. For convenience, we consider a *primed* copy \mathcal{V}' of \mathcal{V} consisting of primed versions of the variables in \mathcal{V} . A *state* is a mapping from \mathcal{V} or \mathcal{V}' to values of the proper type; a *trace* is an infinite sequence of states. We index sequences of states with *instants*, for example, the initial state of a trace is the state at instant 0. Given a formula F over $\mathcal{V} \cup \mathcal{V}'$ and a tuple \mathbf{x} of elements of $\mathcal{V} \cup \mathcal{V}'$, we write $F[\mathbf{x}]$ to indicate that the free variables of F are from \mathbf{x} . Using the standard LTL notion of satisfaction of a formula by a trace, we write $F \models G$ for two LTL formulas F and G over \mathcal{V} to indicate that every trace that satisfies F satisfies G as well. Recall that our model \mathcal{S} is described formally by an LTL formula of the form

$$I[\mathbf{v}] \wedge \text{always } T[\mathbf{v}, \mathbf{v}'] \quad (1)$$

where the elements of \mathbf{v} are from \mathcal{V} . Identifying \mathcal{S} with (1), we consider \mathcal{S} to satisfy a property expressed by a formula F iff the entailment $\mathcal{S} \models F$ holds. We do the same for $\hat{\mathcal{S}}$. We represent *events* of type T as pairs $\langle t, b \rangle$ of values of type T and Bool, with the latter being the Boolean type. An event $\langle t, b \rangle$ is *present* if b is true and *absent* otherwise. We identify any two events of the form $\langle t_1, \text{false} \rangle$ and $\langle t_2, \text{false} \rangle$ since the values t_1 and t_2 are not meaningful.

In our model \mathcal{S} of Uptane, all input and output variables have an event type, to account for the fact that inputs or outputs may be absent. A state of \mathcal{S} has an *eventful input* if there is at least one input event, that is, if the state maps at least one system-level input variable to a value of the form $\langle t, \text{true} \rangle$.

System \mathcal{S} can be described modularly and hierarchically as the parallel composition of several components, with a top-level component standing for the entire system. We can then talk about the input, output and internal variables of each component of \mathcal{S} . Consequently, we say that a state has an *eventful output* for a particular component c of \mathcal{S} if it maps at least one output variable of c to

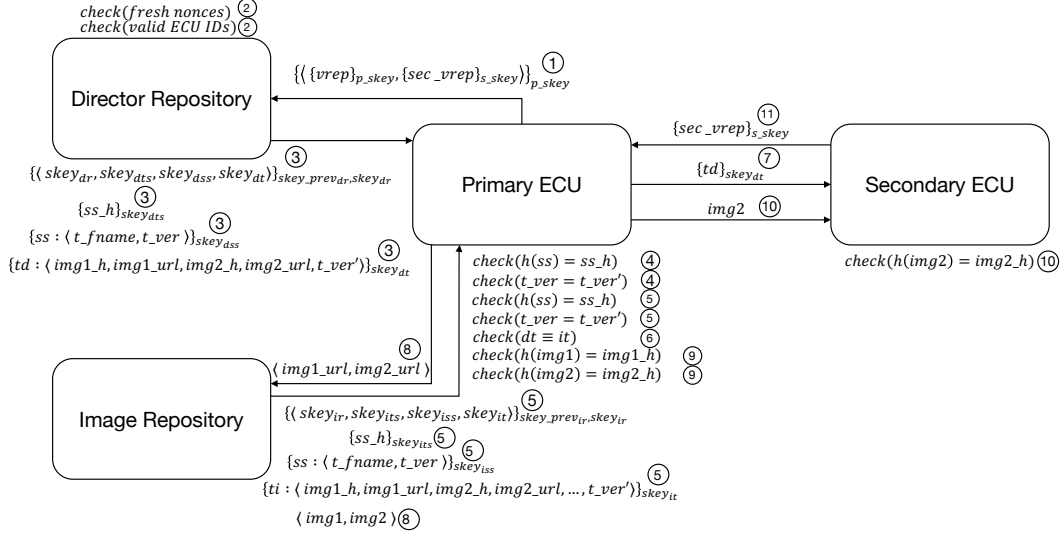


Figure 8: Uptane Alice-Bob Description

a value of the form $\langle t, true \rangle$. We use $Event$ to denote a one-state formula that is satisfied by a trace iff there is an eventful system-level input in the current state. For each system component c we use $Event^c$ to denote a one-state formula that is satisfied by a trace iff there is an eventful output from c in the current state. We will collect all system-level assumptions we make about \mathcal{S} 's input at each state in a formula A .

Note that all properties from Section 3 can be expressed in the form **always** F (see the appendix of the main paper). Since uneventful outputs are not meaningful, we guard with $Event^c$ each property **always** F where F references output from component c and also contains no temporal operators. In other words, we use **always** ($Event^c \Rightarrow F$) instead. We abbreviate $Event^c \Rightarrow F$ as F^c for conciseness and readability.

In the following discussion and in Proposition 1 and Proposition 2, we will focus on models \mathcal{S} and $\hat{\mathcal{S}}$ which express the *benign* case, in the absent of any security attacks. We consider their instrumented versions \mathcal{S}^+ and $\hat{\mathcal{S}}^+$ later in Proposition 3.

In Section D.1, we discuss the specifics of mapping execution traces of $\hat{\mathcal{S}}$ to execution traces in of \mathcal{S} and vice versa. Before that it is useful to recall that in \mathcal{S} , metadata verification occurs across multiple states — eight states are each associated with the verification of a single input metadata file, while the ninth state is used to cross-reference metadata from both repositories. In $\hat{\mathcal{S}}$, all these verification states are collapsed to a single state where verification procedures are applied to all input metadata files simultaneously. We focus on the primary ECU component because the rest of the components are defined in exactly the same way in \mathcal{S} and $\hat{\mathcal{S}}$. To aid in this discussion, Figure 9 and Figure 10 demonstrate the important details of the behavior of the primary ECU component in \mathcal{S} and $\hat{\mathcal{S}}$, respectively. The behavior is presented as a finite state machine where each transition is labeled with (i) a condition on the component's input, (ii) a description of how internal state is updated, and (iii) a description of the component-level outputs (each

separated with a comma). Notice that in both \mathcal{S} and $\hat{\mathcal{S}}$, the primary ECU component is a deterministic FSM— for each set of inputs and current state of the FSM, only a single transition is possible out of that state.

We present a list remarks which are meta-level properties of the models \mathcal{S} and $\hat{\mathcal{S}}$. These remarks will be useful in the development and proof of Proposition 1 and Proposition 2.

REMARK 1. Both \mathcal{S} and $\hat{\mathcal{S}}$ are reactive in the following sense: for all traces of \mathcal{S} (or $\hat{\mathcal{S}}$), a component will produce an eventful output only in states with eventful inputs.

Justification: This is a natural result of Uptane's design. For example, an ECU will not update its current set of metadata (producing an output event) unless it receives some incoming metadata (an eventful input).

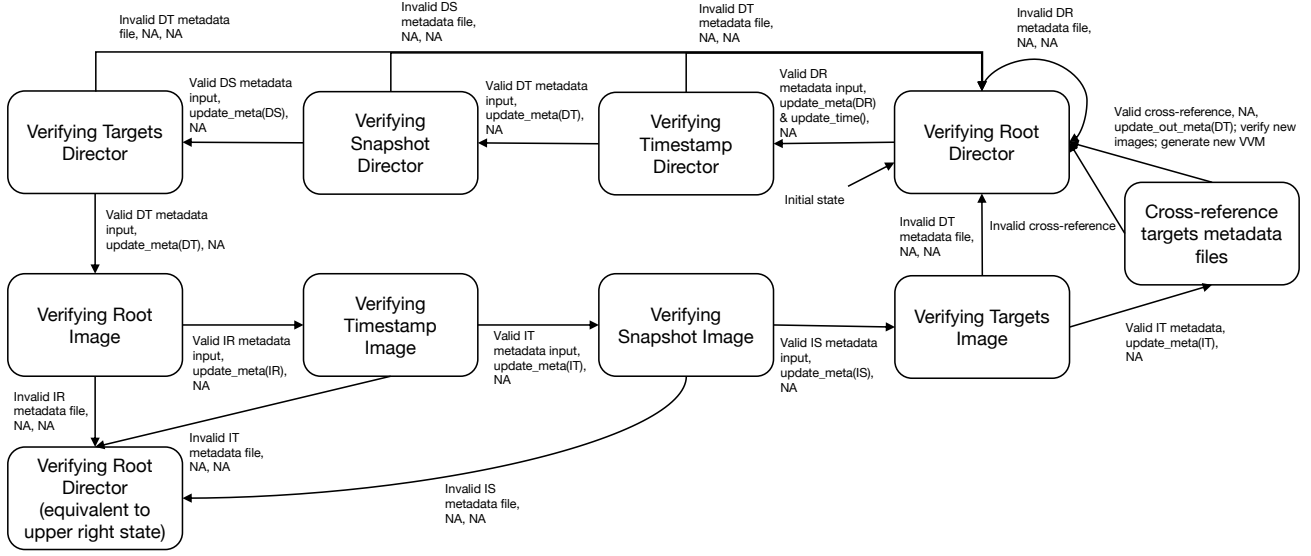
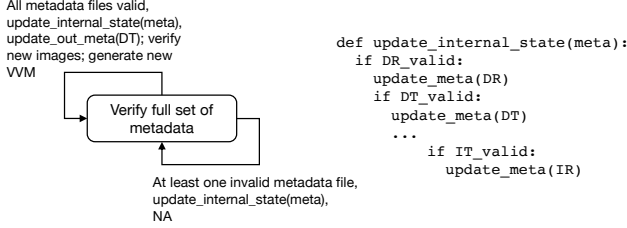
REMARK 2. All ECUs in \mathcal{S} and $\hat{\mathcal{S}}$ are deterministic, that is, given the same input and internal state, they produce the same output.

Justification: This is again by design in Uptane.

REMARK 3. Consider a trace transformation that takes a trace σ of \mathcal{S} and produces a trace $\hat{\sigma}$ by arbitrarily removing states from σ . (A trace transformation like this will be useful when converting a trace in \mathcal{S} to a trace in $\hat{\mathcal{S}}$.) If σ satisfies **always** A where A is a system-level assumption, so will $\hat{\sigma}$.

Justification: All the input assumptions we make on \mathcal{S} based on the Uptane specification are reflexive, transitive relations between consecutive input pairs (e.g., input metadata at state t_i has a version number smaller than or equal to input metadata at state t_{i+1}). Thus, removing states cannot violate these assumptions.

REMARK 4. Consider a trace transformation that takes a trace σ of $\hat{\mathcal{S}}$ and produces trace $\tilde{\sigma}$ by arbitrarily repeating states from σ . (A trace transformation like this will be useful when converting a trace

Figure 9: FSM describing the primary ECU's modes in S .Figure 10: FSM describing primary's modes in \hat{S} .

in \hat{S} to a trace in S .) If σ satisfies always A where A is a system-level assumption, so will $\tilde{\sigma}$.

Justification: As with Remark 3.

D.1 Converting traces between S and \hat{S}

In this section, we describe how to convert traces of the uncompressed system S to traces of the compressed system \hat{S} and vice versa, which is necessary to prove Proposition 1.

For simplicity, we make the following assumption, with possible loss of generality:

ASSUMPTION 1. *At any one time, either all system-level input is present or none is.*

D.1.1 Trace expansion. We describe how to expand a trace σ of \hat{S} to a trace $\tilde{\sigma}$ of S . Trace expansion is necessary for the proof of Proposition 1. We proceed by replacing each state s in σ that represent a (compressed) verification cycle with a sequence of states s representing the uncompressed version of that cycle. First, assume the verification cycle succeeded in s . Notice that the primary in the uncompressed model must have at least nine states in each complete verification cycle, and that each state processes a separate input

metadata file. So, to construct the input for s , we must replicate the input from s across nine states. This way, at each state, the primary ECU has still access to the correct metadata file in its input. Since the input is fixed and the uncompressed primary ECU is deterministic, we can compute the value of the internal and output variables of each state in s based on Figure 9. Notice that the uncompressed primary ECU will produce uneventful outputs until the last state in the verification cycle, that is, the transition from the “Cross-referencing” state to the “Verifying Root Director” state. In that last state, the uncompressed primary will output its updated metadata and images done by the compressed primary ECU in state s (where metadata and images are instantaneously updated). Also, the uncompressed primary ECU will update its internal state at every state of the cycle, ultimately resulting in the primary ECU completely updating its current set of metadata (which also happens in the compressed primary in s).

If verification failed in s , then to construct the input for s we only need to replicate the input n times, where n is the number of metadata files processed by the uncompressed primary ECU before one of them fails verification. We need n states in the uncompressed model (Figure 9) because the uncompressed primary ECU has a state for every metadata file until either (i) every metadata file is verified, or (ii) a metadata file fails verification, returning the uncompressed primary ECU to the “Verifying Director Root” state. In case (ii), this results in n states for the verification cycle where no outputs are produced (corresponding to the non-event output produced by the compressed primary in s). The uncompressed primary only updates its internal state for the first $n - 1$ states. Notice that this change in internal state is analogous to the change in the compressed primary, as the `update_internal_state` function will update the internal state for each metadata file up to the file that fails verification. This means that the ECU's internal state will be *analogous* in both models when beginning the next update cycle, where by “analogous” we mean the following.

One of the primary ECU's component-level inputs is a clock value representing the current time. In both \mathcal{S} and $\widehat{\mathcal{S}}$, the clock value is modeled as an internal integer variable that is strictly increasing in each state transition (although, for generality, not by any set amount). The primary ECU processes this clock value input at the beginning of each verification cycle in order to determine whether metadata files (which have expiration times) have expired. When replacing a verification state s from a trace in $\widehat{\mathcal{S}}$ with a verification cycle s , we slightly alter each system-level input metadata file by multiplying its expiration time by nine (in every state of s). Also, we correspondingly multiply the clock value input by nine in the first state of s , while setting the clock to increase by one in each subsequent state of s . For example, if a trace in the compressed model has a clock value of 1 in the first state and a clock value of 2 in the second state, then the corresponding trace in the uncompressed model has a clock value of 9 in the first state of the first uncompressed verification cycle and a clock value of 18 in the first state of the second uncompressed verification cycle. This is necessary because a direct translation of the clock value without multiplication by nine would require the clock value to stutter (keep the same value) during the verification cycle s , which is not allowed in \mathcal{S} .

Observe now that the trace produced in the expansion process just described has a sequence of valid inputs by Remark 4. This leads to the following result.

LEMMA 1. *The trace generated by the expansion process from a trace of $\widehat{\mathcal{S}}$ is a trace of \mathcal{S} .*

D.1.2 Trace compression. We discuss an inverse process that translates a trace σ in \mathcal{S} to a corresponding trace $\widehat{\sigma}$ in $\widehat{\mathcal{S}}$. This sort of trace compression is useful in the proof of Proposition 1. We proceed by replacing each sequence s of states representing a full verification cycle in σ with a single state s in $\widehat{\sigma}$.

First, assume the verification cycle s is successful. Notice that in the uncompressed model, we verify a single metadata file at each state. So, to construct the input for s , we take the metadata file that is being verified in the corresponding state of s , as well as the clock value from the first state in s . Also notice that input images are only processed in the last state of the verification cycle. So, we take the input images from the last state of s . With these inputs, the compressed primary ECU will complete a successful verification cycle and instantaneously produce its updated set of metadata and images as outputs (which happens in the uncompressed primary only in the last state of s). Additionally, it will instantaneously and completely update its internal state, matching the uncompressed primary's internal state at the end of s .

If s is unsuccessful, then we use the same input for s . The only exception to this is if s is unsuccessful because the primary tried to verify a metadata file that was absent. In this case, we can use any metadata file that will fail the primary's verification. With this input, verification will fail and the compressed primary ECU will produce an uneventful output (corresponding to the uneventful outputs produced by the uncompressed primary at every state in s). Additionally, the compressed primary ECU will update its internal state for the first $n - 1$ metadata files (assuming the n -th file fails verification), which is analogous to the updates of internal state

performed by the uncompressed primary across the first $n - 1$ states of s .

By Remark 3 we can claim that every trace produced by the compression process has a sequence of valid inputs.

LEMMA 2. *The trace generated by the compression process from a trace of \mathcal{S} is a trace of $\widehat{\mathcal{S}}$.*

A clarification on the two translation processes is in order. Notice that a state s is the result of compressing a verification cycle s or the sequence s is the result of expanding a verification state s , then each state in s either produces uneventful output (before the last state of s) or produces output (at the last state of s) analogous to the output in s . More precisely, the two output are the same modulo their metadata expiration times (which are altered in the expansion process). That difference is immaterial because the desired system-level properties do not reference metadata expiration times. Moreover, system-level outputs directly correspond to the output of the primary ECU component. This is because system-level outputs can only be produced when the primary is producing component-level outputs, and the remaining components of the model have the exact same behavior in both the compressed and the uncompressed case.

D.2 Proof of Main Theorem

First, we prove that \mathcal{S} and $\widehat{\mathcal{S}}$ are equivalent with respect to the properties under analysis (model compression).

PROPOSITION 1. *For each desired functional requirement P always from Section 3, $\mathcal{S} \wedge \text{always } A \models \text{always } P$ if and only if $\widehat{\mathcal{S}} \wedge \text{always } A \models \text{always } P$.*

PROOF. It is enough to consider traces that satisfy the system-level input assumptions A .

We have two cases depending on the form of $\text{always } P$.

Case one: P is a state property φ^c . Assume φ^c only references component-level output directly from c or from components that directly provide input to c . (If a component c_1 provides input to c_2 , then an output event at c_2 necessarily implies an output event at c_1 due to Remark 1. Hence, if φ^c references output variables from both components, then Event^{c_2} is enough to know that neither component is currently producing an uneventful output.)

\Rightarrow : Assume $\mathcal{S} \wedge \text{always } A \models \text{always } \varphi^c$. Take an arbitrary trace σ in $\widehat{\mathcal{S}}$ that satisfies input assumptions. First, we construct a sequence of states $\check{\sigma}$ by expanding each state to a verification cycle across multiple states as described in Section D.1. We know that $\check{\sigma}$ is a trace of \mathcal{S} by Lemma 1.

Since $\check{\sigma}$ is a trace of \mathcal{S} , φ^c must hold in every state of $\check{\sigma}$. Also, each state in σ has equivalent output to some state in $\check{\sigma}$, and φ^c has no temporal operators, so φ^c must hold for all states of σ . Therefore, $\widehat{\mathcal{S}} \wedge \text{always } A \models \text{always } \varphi^c$.

\Leftarrow : Assume $\widehat{\mathcal{S}} \wedge \text{always } A \models \text{always } \varphi^c$. Take an arbitrary trace σ in \mathcal{S} . We will construct a sequence of states $\hat{\sigma}$. First, collapse each verification cycle in σ into a single state as described in Section D.1. We know that $\hat{\sigma}$ is a trace of $\widehat{\mathcal{S}}$ satisfying system-level input assumptions by Lemma 2.

Since $\hat{\sigma}$ is a trace of $\widehat{\mathcal{S}}$, P must hold at every state of $\hat{\sigma}$.

Now, we show σ satisfies φ^c at every state. First, we know that φ^c is satisfied in every state present in $\hat{\sigma}$. All other states have non-event outputs, so they trivially satisfy φ^c . Therefore, $\mathcal{S} \wedge \text{always } A \models \text{always } \varphi^c$.

Case two: P is a property of the form **always** (historically $\varphi_1^c \Rightarrow \varphi_2^c$). Assume φ_1^c and φ_2^c only reference component-level output directly from c .

\Rightarrow : Assume $\mathcal{S} \wedge \text{always } A \models \text{always (historically } \varphi_1^c \Rightarrow \varphi_2^c)$. Take an arbitrary trace σ in $\hat{\mathcal{S}}$. Let t_i be the first state in σ where φ_1^c evaluates to **false**. (If φ_1^c never evaluates to **false**, then P is equivalent to **always** φ_2^c , falling under case one.) First, we construct a sequence of states $\hat{\sigma}$ by expanding each verification cycle as described in Section D.1. We know that $\hat{\sigma}$ is a trace of \mathcal{S} by Lemma 1.

Let t_j be the state in $\hat{\sigma}$ with output equivalent to t_i in σ , and notice that φ_1^c cannot be falsified in $\hat{\sigma}$ before t_j . (If φ_1^c was falsified before t_j , then there must be a corresponding state in σ that falsifies φ_1^c before t_i , a contradiction.) Since $\hat{\sigma}$ is a trace of \mathcal{S} , it must satisfy **always** P .

Now, we show σ satisfies P at every state. Each state in σ at or after t_i trivially satisfies **historically** $\varphi_1^c \Rightarrow \varphi_2^c$, since t_i falsifies φ_1^c . Also, each state in σ before t_i either has an uneventful output or has equivalent output to a state in $\hat{\sigma}$ before t_j , which must satisfy φ_2^c . Therefore, $\hat{\mathcal{S}} \wedge \text{always } A \models \text{always } P$.

\Leftarrow : Assume

$$\hat{\mathcal{S}} \wedge \text{always } A \models \text{always (historically } \varphi_1^c \Rightarrow \varphi_2^c).$$

Take an arbitrary trace σ in \mathcal{S} , and let t_i be the state in σ where φ_1^c first evaluates to **false**. (If φ_1^c never evaluates to **false**, then P is equivalent to **always** φ_2^c , falling under case one.) We will construct a sequence of states $\hat{\sigma}$. First, collapse each verification cycle in σ into a single state as described in Section D.1. We know that $\hat{\sigma}$ is a trace of $\hat{\mathcal{S}}$ satisfying system-level input assumptions by Lemma 2.

Let t_j be the state in $\hat{\sigma}$ with output equivalent to t_i in σ , and notice that φ_1^c cannot be falsified in $\hat{\sigma}$ before t_j . Since $\hat{\sigma}$ is a trace of $\hat{\mathcal{S}}$, it must satisfy **always** P .

Now, we show σ satisfies P at every state. First, every state in σ after t_i trivially satisfies **historically** $\varphi_1^c \Rightarrow \varphi_2^c$, since t_i falsifies φ_1^c . Also, every state in σ before t_i is either an uneventful output or has equivalent output to a state in $\hat{\sigma}$ before t_j , which must satisfy φ_2^c .

Therefore,

$$\mathcal{S} \wedge \text{always } A \models \text{always (historically } \varphi_1^c \Rightarrow \varphi_2^c).$$

P1, P2, P3, P4, P5, P6, and P8 fall under case one. P7 falls under case two. \square

Next, we prove that in $\hat{\mathcal{S}}$, with respect to the properties under analysis, we can assume that there is an event input at every state without loss of generality (input compression).

PROPOSITION 2. *For each desired functional requirement **always** P from Section 3, $\hat{\mathcal{S}} \wedge \text{always } A \models \text{always } P$ if and only if $\hat{\mathcal{S}} \wedge \text{always } (A \wedge \text{Event}) \models \text{always } P$.*

PROOF. \Rightarrow : Since $\hat{\mathcal{S}} \wedge \text{always } (A \wedge \text{Event}) \models \text{always } P$ is a strictly weaker claim, the result immediately follows.

\Leftarrow : We will only consider traces that satisfy the system-level input assumptions A .

We have two cases depending on the form of **always** P .

Case one: P is a state property φ^c . Assume φ^c only references component-level output directly from c or from components that directly provide input to c . (If a component c_1 provides input to c_2 , then an output event at c_2 necessarily implies an output event at c_1 due to Remark 1. Hence, if φ^c references output variables from both components, then Event^{c_2} is enough to know that neither component is currently producing an uneventful output.)

Assume $\hat{\mathcal{S}} \wedge \text{always } (A \wedge \text{Event}) \models \text{always } \varphi^c$. Take an arbitrary trace σ in $\hat{\mathcal{S}}$. We will construct a sequence of states $\hat{\sigma}$ by removing states with non-event inputs. This leaves a valid trace of $\hat{\mathcal{S}}$ by Remark 2 and Remark 3.

Since $\hat{\sigma}$ is a trace of $\hat{\mathcal{S}}$ and has event inputs at every state, P must hold at every state of $\hat{\sigma}$.

Now, we show σ satisfies φ^c at every state. First, we know that φ^c is satisfied in every state present in $\hat{\sigma}$. All other states must have non-event outputs (by Remark 1), so they trivially satisfy φ^c . Therefore, $\mathcal{S} \wedge \text{always } A \models \text{always } \varphi^c$.

Case two: P is a property of the form **always** (historically $\varphi_1^c \Rightarrow \varphi_2^c$). Assume φ_1^c and φ_2^c only reference component-level output directly from c .

Assume

$$\hat{\mathcal{S}} \wedge \text{always } (A \wedge \text{Event}) \models \text{always (historically } \varphi_1^c \Rightarrow \varphi_2^c).$$

Take an arbitrary trace σ in $\hat{\mathcal{S}}$, and let t_i be the state in σ where φ_1^c first evaluates to **false**. (If φ_1^c never evaluates to **false**, then P is equal to **always** φ_2^c , falling under case one.) We will construct a sequence of states $\hat{\sigma}$ by removing states with non-event inputs. This leaves a valid trace of $\hat{\mathcal{S}}$ by Remark 2 and Remark 3.

Let t_j be the state in $\hat{\sigma}$ with output equal to t_i in σ , and notice that φ_1^c cannot be falsified in $\hat{\sigma}$ before t_j . Since $\hat{\sigma}$ is a trace of $\hat{\mathcal{S}}$ and has event inputs at every state, it must satisfy **always** P .

Now, we show σ satisfies P at every state. First, every state in σ after t_i trivially satisfies **historically** $\varphi_1^c \Rightarrow \varphi_2^c$, since t_i falsifies φ_1^c . Also, every state in σ before t_i is either an uneventful output (by Remark 1) or has equal output to a state in $\hat{\sigma}$ before t_j , which must satisfy φ_2^c .

Therefore, $\mathcal{S} \wedge \text{always } A \models \text{always (historically } \varphi_1^c \Rightarrow \varphi_2^c)$.

P1, P2, P3, P4, P5, P6, and P8 fall under case one. P7 falls under case two. \square

Recall that we denote the instrumented version of \mathcal{S} as \mathcal{S}^+ , which incorporates adversarial influence into the model. Similarly, we denote the instrumented version of $\hat{\mathcal{S}}$ as $\hat{\mathcal{S}}^+$. In the instrumented models, the adversary is capable of modifying the value of data that is sent between components. The following proposition demonstrates that property violations in the compressed model correspond to property violations in the uncompressed model.

PROPOSITION 3. *For each desired functional requirement R from Section 3, $\hat{\mathcal{S}}^+ \wedge \text{always } (A \wedge \text{Event}) \models R \Rightarrow \mathcal{S}^+ \wedge \text{always } A \models R$*

PROOF. Assume $\hat{\mathcal{S}}^+ \wedge \text{always } (A \wedge \text{Event}) \models R$. Then there exists some trace σ in $\hat{\mathcal{S}}^+$ with eventful inputs at every state which falsifies R . We will argue that there exists a corresponding trace $\tilde{\sigma}$ of \mathcal{S}^+ that also falsifies R . The general idea is to take σ and expand it according to Lemma 1, but to also include a corresponding adversarial modification for each adversarial action applied in σ .

Suppose in σ that the adversary modified the value of a variable sent between component A and component B in state s . Then, in our expanded trace $\tilde{\sigma}$, we replace s with a verification cycle s consisting of a sequence of states. If component A is a primary ECU and component B is a secondary ECU, then the corresponding adversarial action is applied at the last state of s (the state in which messages are exchanged, as the ECUs do not produce output during the verification cycle). If component A is either repository and component B is the primary ECU, then similarly, the adversarial action is applied at the last state of s , *unless* the adversarial action was a modification of a metadata file sent from a repository to the primary. If that is the case, then the adversarial action is applied to the state where the primary's verification mode matches the metadata file being modified. For example, if state s in σ includes an adversarial action modifying the director repository's snapshot metadata, then the action is applied in the third state of s in $\tilde{\sigma}$.

Now, we argue that $\tilde{\sigma}$ is indeed a valid trace of S^+ . First, the sequence of inputs is a valid sequence of inputs by Lemma 1. Second, every adversarial action in \tilde{S}^+ is a possible adversarial action in S^+ , as the adversary's capabilities are identical in both models. Third, $\tilde{\sigma}$ contains a valid sequence of outputs with respect to the behavior of S^+ . If an ECU in the compressed model processes some incoming metadata file m in σ at state s (where m is potentially injected by an adversary), we copy the adversarial action such that the corresponding state s' in the corresponding verification cycle s in $\tilde{\sigma}$ is also processing m (modulo metadata expiration times). From Remark 2 we know that ECUs are deterministic, so the ECU processing m at state s behaves identically to the ECU processing m at state s' (that is, either m is verified in both situations or m is rejected in both situations). Thus, the outputs in $\tilde{\sigma}$ must be valid with respect to S^+ .

With the above adversarial actions incorporated into $\tilde{\sigma}$, every state in σ must have a state with equal output in $\tilde{\sigma}$ (modulo metadata expiration times). We have two cases for the form of R :

- (1) R is of the form **always** P , where P is a state property. In this case, there is a property violation of σ at state t_i . Because $\tilde{\sigma}$ contains some state with equivalent output t_j (modulo metadata expiration times), and because P is a state property, $\tilde{\sigma}$ is violated at t_j .
- (2) R is of the form **always** (historically $\varphi_1^c \Rightarrow \varphi_2^c$), where φ_1^c and φ_2^c are state properties. A property violation of R in σ means that φ_2^c was false in some state t_i before φ_1^c ever evaluated to false at a later state t_j (if at all). Then there must be a corresponding state t_k in $\tilde{\sigma}$ that falsifies φ_2^c . Further, $\tilde{\sigma}$ could not have falsified φ_1^c before t_k , or else there would be a corresponding falsification of φ_1^c before t_i in σ .

□

In the other direction, we show that property violations in the uncompressed model correspond to property violations in the compressed model.

PROPOSITION 4. *For each desired functional requirement R from Section 3, $S^+ \wedge \text{always } A \models R \Rightarrow \hat{S}^+ \wedge \text{always } (A \wedge \text{Event}) \models R$*

PROOF. Assume $S^+ \wedge \text{always } A \not\models R$. Then there exists some trace σ in S^+ which falsifies R . We will argue that there exists a corresponding trace $\hat{\sigma}$ of \hat{S}^+ with eventful inputs at every state that

also falsifies R . The general idea is to take σ and compress it, but to also include a corresponding adversarial modification for each adversarial action applied in σ .

Instead of using Lemma 2 to compress σ , we will perform a more straightforward conversion. The main idea is to produce a trace $\hat{\sigma}$ such that \hat{S}^+ simulates S^+ , *i.e.*, each n -state verification cycle in σ maps to an n -state "pseudo"-verification cycle in $\hat{\sigma}$ where only one metadata file is updated at a time, even in the compressed model. We now describe the conversion.

Consider an arbitrary state s in σ which corresponds to some instant i in a verification cycle. We will replace s with s' , the corresponding state in $\hat{\sigma}$. In fact, s' is equivalent to s , except for the following case. For communications between the repositories and the primary ECU, include the following adversarial actions:

- (1) For metadata files at instants 0 through $i - 1$, replay files that are already present in the ECU's state.
- (2) For the metadata file at instant i , copy the adversarial action from s .
- (3) For all metadata files at all other instants greater than i , inject garbage values that will fail the ECU's verification.

We argue that the above process indeed produces a trace of \hat{S}^+ . First, the sequence of inputs is unchanged and is therefore still valid. Second, the above adversarial actions are possible no matter the specific type of adversarial instrumentation— the adversary either performs a replay, copies the adversarial action from σ , or injects garbage values. Third, we argue that the sequence of outputs is valid. This holds because the effect of the above actions is that the ECU in $\hat{\sigma}$ will only (potentially) update one metadata file per state, which corresponds exactly to the metadata file updated in the corresponding state of σ .

Note that for action (1) (replaying earlier metadata files), metadata expiration times may cause files to fail verification. However, this issue can be sidestepped by increasing the expiration time of the metadata file, causing it to pass verification but with no effect on the system-level properties which do not directly reference metadata expiration times.

Recall also that in σ , there may be states without eventful input, but that we want to produce a trace $\hat{\sigma}$ with eventful input at every state. Similarly to the previous strategy, we can simulate the absence of eventful input by having the adversary inject garbage values in every communication channel, which causes the same effect on ECU state.

With the above adversarial actions incorporated into $\hat{\sigma}$, the values of the output variables are equivalent in σ and $\hat{\sigma}$. Therefore, a property violation in σ directly results in a property violation in $\hat{\sigma}$. □

We are now ready to prove our main theorem.

THEOREM 2. *Let A be a formula denoting assumptions on system-level inputs for S . For each desired functional requirement R from Section 3:*

- (a) $S \wedge \text{always } A \models R$ iff $\hat{S} \wedge \text{always } (A \wedge \text{Event}) \models R$
- (b) $\hat{S}^+ \wedge \text{always } (A \wedge \text{Event}) \models R$ iff $S^+ \wedge \text{always } A \models R$

PROOF. (a) Recall that R has the form **always** P where P is a state formula. From Proposition 1 we have $S \wedge \text{always } A \models \text{always } P$ if and only if $\hat{S} \wedge \text{always } A \models \text{always } P$, and from Proposition 2

we have $\hat{S} \wedge \text{always } A \models \text{always } P$ if and only if $\hat{S} \wedge \text{always } (A \wedge \text{Event}) \models \text{always } P$. Composing the two bi-implications, we get $S \wedge \text{always } A \models \text{always } P$ if and only if $\hat{S} \wedge \text{always } (A \wedge \text{Event}) \models \text{always } P$.

(b) We have $\hat{S}^+ \wedge \text{always } (A \wedge \text{Event}) \not\models R \Rightarrow S^+ \wedge \text{always } A \not\models R$ from Proposition 3, and we have $S^+ \wedge \text{always } A \not\models R \Rightarrow \hat{S}^+ \wedge \text{always } (A \wedge \text{Event}) \not\models R$ from Proposition 4. \square

E IMPLEMENTATION

In this section, we describe the details of our approach.

E.1 Modeling Adversarial Influence

To model adversarial influence, we place an adversarial component A between each pair of components (M, N) in S . When M sends a message to N (or vice versa), A takes in M 's output and nondeterministically modifies it before passing it along to N , acting as a man in the middle.

The adversary's behavior is defined as an abstract description of its output values in terms of its input values, called a *contract*. This allows for nondeterminism, where the adversary is free to pick any possible modification that respects the contract.

Standard Attacker. No restrictions are imposed on the output value of the component, so the adversary can select an arbitrary message to inject.

Bounded Replay Attacker. The bounded replay attacker can only inject messages that were sent along the network connection in the last k timesteps (for some constant k). The adversary's contract involves storing each message sent from M to N in a wraparound queue of size k . This wraparound queue can be seen as a sliding window where the attacker only holds the last k messages in its memory. Then, the attacker must pick an element from this window when injecting messages along the network channel. This disallows the adversary from *arbitrarily* modifying the data sent from A to B , as the adversary can now only replay one of the last k messages.

Unbounded Replay Attacker. We also model a replay attacker with unbounded memory, *i.e.*, the ability to replay messages that were sent arbitrarily far in the past. Instead of storing previous messages in a queue, the unbounded replay attacker stores the values of previous messages as outputs of a partial function $f : \mathbb{N} \rightarrow B$, where B is the type of the message being sent along the channel. At each timestep i , the adversary extends the definition of f such that $i \mapsto \text{msg}_i$, where msg_i is the current message being sent along the channel. Then, at timestep i , the adversary is free to nondeterministically select any message $m \in \{f(j) \mid 0 \leq j \leq i\}$.

E.2 MC Model

Uptane System Architecture. We begin by outlining Uptane's system architecture which is graphically demonstrated in Figure 2. We define the top level component and each of its subcomponents and then specify each component's interface and connections. For example, here is the SecondaryECU component's interface:

```
component SecondaryECU
  (in_primary: PrimaryToSecondary)
  returns
    (out_primary: SecondaryToPrimary,
     installed_image_secondary: Image,
```

```
     verified_metadata_secondary: Metadata)
```

The SecondaryECU component has one input where it receives data from the primary ECU and three outputs where it (i) sends data to the primary ECU, (ii) reports its currently installed image, and (iii) reports its current set of verified metadata. The `in_primary` and `out_primary` variables contain the data being sent from the primary to the secondary, and vice versa (*e.g.* metadata, ECU version reports). The `installed_image_secondary` and `verified_metadata_secondary` variables are also records, the former containing information about an image and the latter containing four metadata files.

Component-Level Design. Next, we specify the behavior of each component of the architecture with assume-guarantee contracts. Consider the following functional specifications for the SecondaryECU: (1) The ECU's version report has correct information with respect to the currently installed image (filename and hashes); and (2) the ECU updates metadata according to the Uptane standards document.

We formalize these specifications using abstract syntax for conciseness. First, for (1), we add the following guarantee stating that the filename and hash listed in the output ECU version report equal the installed image's filename and hash (`image` abbreviates `installed_image_secondary`):

```
image.filename = out_primary.report.filename and
hash(image) = out_primary.report.hash
```

To formalize aspect (2), we add the following guarantee where `image` abbreviates `installed_image_secondary`:

```
if new_image_verified
then image' = new_image
else image' = image
```

where `image'` denotes the value of variable `image` in the next state, and `new_image_verified` is a predicate defined according to the image verification instructions in the Uptane standards document. The guarantee states that if the new image passes verification, then the ECU updates to the new image; otherwise, it keeps the previously-installed image. In the above example, the initial state predicate enforces that `image` is equal to `initial_image` initially.

Any system execution that satisfies these constraints will be considered valid during the analysis performed by MC. We follow a similar process to specify the other components.

Desired Functional Requirements. Recall the list of desired functional requirements from Section 3. We can express, for example, property P3 (checks for rollback attacks) with the following top-level guarantee, where the primed version of `installed_image_secondary` refers to the value of the variable in the next state:

```
installed_image_secondary'.version >=
installed_image_secondary.version
```

Simplifying Assumptions. To make the analysis feasible, we make further simplifying assumptions: (i) We assume there are no delegations¹ in targets metadata; (ii) we model a vehicle with exactly one primary and two secondary ECUs; and (iii) we omit parts of the standard that relate to implementation details, *e.g.*, how filenames are encoded.

¹Delegations refer to when signing ability is deferred to another party. See the Uptane standard [27] for more details. To the best of our knowledge, support for delegations is not yet present in open source implementations.

E.3 CPV Model

To accurately model the adversary's capabilities with respect to cryptographic assumptions, we consult Tamarin [50], a cryptographic protocol verifier. We formulate a Tamarin model by including rules capturing the cryptographic aspects of each message sent in the Uptane protocol (*i.e.*, each connection in the MC model). As an example, the following rule models the primary ECU sending targets metadata to the secondary ECU:

```
rule primary_send_metadata:
  [!PriLatestTargetsDir(targets, <sig1, sig2>)]
  -->
  [Out(< targets, sig1, sig2 >)]
```

The primary ECU retrieves the most recently verified targets metadata (and its signatures) and sends it along the insecure network channel. The secondary ECU verifies metadata with the following rule:

```
rule secondary_verify_metadata:
  [ In(<targets, sig1, sig2>),
    SecMetaDir(<pubkey1, pubkey2>, old_targets) ]
  --[ Eq(verify(sig1, targets, pubkey1), true),
        Eq(verify(sig2, targets, pubkey2), true), ]->
  [SecMetaDir(<pubkey1, pubkey2>, targets)]
```

The secondary ECU receives the incoming metadata from the insecure channel and performs verification by checking the digital signatures. The Tamarin model only captures the *cryptographic aspects* of the protocol while abstracting away other aspects of verification that are present in the model \mathcal{S} (*e.g.*, checking version numbers).

In the Tamarin model, we formulate 10 correspondence (weak authentication) and 10 injective-correspondence (strong authentication) lemmas to learn about the cryptographic aspects of the protocol. The correspondence lemmas are of the form:

$$\forall msg. \forall i. \text{Receive}(msg)@i \Rightarrow (\exists j. \text{Send}(msg)@j \wedge (i < j)).$$

In other words, if a message (*e.g.*, metadata) was received (*i.e.*, verified) at time i , then it must have been sent at an earlier time j . If this lemma is proven for some network connection, then it is cryptographically infeasible for the adversary to inject new messages.

The injective-correspondence lemmas are of the form:

$$\forall m. \forall i. \text{Receive}(m)@i \Rightarrow (\exists j. \text{Send}(m)@j \wedge (j < i) \wedge \neg(\exists i2. \text{Receive}(m)@i2 \wedge \neg(i2 = i))).$$

This lemma states that if a message was received at time i , then it must have a *unique* matching sender at an earlier time j . If this lemma is proven, then the adversary cannot inject *or* replay messages; failure suggests that replays are possible. The two previous lemmas are formulated for each step of metadata verification, as well as for VVM and image verification. In addition to the authentication lemmas, it is necessary to prove *sanity-check* lemmas of the form $\exists msg. \exists i. \text{Receive}(msg)@i$. that demonstrate that the premises of the implications are reachable. These lemmas are called a “sanity-check” lemma because if it is disproven by Tamarin, then there is likely a mistake in the model. We proved sanity-check lemmas for every protocol connection.

Termination. Tamarin often has to find long traces to (dis)prove lemmas, resulting in termination issues. We thus implemented our own heuristic (called an *oracle*) to guide Tamarin in its proof search.

This heuristic orders Tamarin's proof goals, instructing Tamarin on which unsolved premises to solve first. Our heuristic is Uptane-specific and sets proof goals in the logical order that metadata verification occurs in.

Another essential step in achieving termination is to impose bounds on the number of times each Tamarin rule can be applied (in a trace). In some cases, for tractability, we only allow each Tamarin rule to apply only once or twice in a trace. Although this is limiting, Tamarin was still able to disprove several lemmas, leading to the attacks we found.

F FINDINGS

F.1 Working Example, Finding 2 and Finding 3 (VVM Replays)

To further illustrate the workflow, we will walk through the analysis of desired functional requirement P2. In English, the property we want to check is “The director repository never verifies old VVMs.” In the attack, the adversary intercepts and replays an old VVM, tampering with the director repository's task of generating new metadata (potentially causing it to direct the vehicle to install the wrong images).

Execute CPV. First, we execute the relevant lemmas in the Tamarin model— those corresponding to (i) sending the VVM from the primary ECU to the director, and (ii) sending an ECU version report from the secondary to the primary. In both cases, the weak authentication lemma is proven, meaning that injection of new VVMs/version reports is cryptographically infeasible. Also, in both cases, the strong authentication lemma fails with counterexample traces where the adversary replays the same VVM/ECU version report twice, and they are still verified by the director.

Specify MC threat model. The results from the CPV model (weak authentication lemma proven, strong authentication lemma disproven) mean that a *replay attacker* should be placed on the connections from the primary ECU to the director and from the secondary ECU to the primary ECU. All other connections are labeled as invulnerable to attack.

Execute MC. MC produces a counterexample where an adversary intercepts a VVM and replays it at a later timestep, a straightforward replay that represents a violation of P6.

When the director receives the VVM, it checks that the VVM is valid by verifying the digital signatures, comparing the ECU IDs to its internal database, and checking if the nonces are fresh (based on if they have been seen before). However, the nonce check is not effective against an adversary who intercepts a VVM and replays it at a later time. The check still succeeds, as this is still the first time the director has seen the nonce (assuming the adversary does not replay the same VVM multiple times). Since the replayed VVM still contains valid signatures and ECU IDs, the replay succeeds.

Attack chaining. We perform a technique called *attack chaining* where we update the model to block the previous attack trace and see if the tool can find a different (and potentially more interesting) attack. We update the model by (i) updating the system-level property to “The director never verifies an old ECU version report,” a more fine-grained property, and by (ii) modeling the connection from the primary ECU to the director as invulnerable to attack. These changes guarantee a trace that is different than the previous

attack. With these updates to the model, MC produces a counterexample trace representing a unique attack. In this version, the adversary dropped an ECU version report being sent from a secondary to the primary ECU and injected it at a later cycle. In this version of the attack, the adversary created a VVM from ECU version reports that weren't initially supposed to be in the same VVM.

When the primary builds the VVM at the beginning of each update cycle, it queries all its secondaries for their latest ECU version reports. However, the primary does not do any verification (all verification is performed by the director). The director thus is still vulnerable to replay attacks. Instead of replaying at the manifest level, the adversary can intercept and replay individual version reports within the vehicle. When the ECU version reports are replayed, the primary still signs the manifest, so all the signatures are still valid.

F.2 Other Findings

Finding 4: Arbitrary Software Installation Attack. The adversary causes a secondary ECU to install an adversary-authored software image.

Threat Model. Internal adversary with compromised director targets metadata keys (A3).

Vulnerability. When the secondary ECU performs partial verification on metadata, there is no cross-referencing of metadata from the image repository. The compromise of only director targets metadata keys (no image repository keys) leaves the secondary ECU vulnerable to any attack. If the adversary performed a man-in-the-middle attack solely outside the vehicle, the primary would detect the attack (when cross-referencing metadata from both repositories) and would refrain from forwarding the offending targets metadata file to any secondary ECU. But, an internal adversary can directly inject a malicious metadata file, bypassing the primary ECU's verification.

Detection. Based on Tamarin's output, we insert a standard attacker on the connection from the primary to the secondary. With the standard attacker, MC finds a counterexample to the property "The secondary ECU never verifies an attacker-authored image." In the counterexample, the adversary first forges and injects a targets metadata file to the secondary ECU that contains the hash of an attacker-authored image. The forged metadata file is verified by the ECU because the digital signatures are valid and the adversary has complete control over the metadata file's contents. Then, the attacker injects the corresponding attacker-authored image, which passes the secondary's verification because its hash matches the hash found in the secondary's (forged) targets metadata.

Finding 5: Incompatible Image Installation Attack. The adversary causes an ECU to install an image that is incompatible with the ECU's hardware. This is different from the incompatibilities discussed in previous attacks, which were in the form of images being incompatible with *each other*, not with the ECU's hardware.

Threat Model. Internal adversary with compromised director targets metadata keys (A3).

Vulnerability. Same as the arbitrary software attack.

Detection. Similar to arbitrary software attack.

Finding 6: Mix-and-Match, Rollback Attacks. We also analyzed mix-and-match attacks and rollback attacks. These attacks are possible in attack scenario A3 due to the same vulnerability that

Attack Type	Possible Implications
Freeze (P1)	Expose deprecated vulnerabilities
VVM Replay (P2)	Expose deprecated vulnerabilities, hampered vehicle functionality
Rollback (P3)	Expose deprecated vulnerabilities
Arbitrary Software (P4)	Adversary control over vehicle
Attacker-Authored VVM (P5)	Expose deprecated vulnerabilities
Mix-And-Match (P6)	Hampered vehicle functionality
Mixed-Bundles (P7)	Hampered vehicle functionality
Incompatible Image (P8)	Hampered vehicle functionality

Table 8: Possible implications of attacks

allows arbitrary software attacks. But, they have a lower impact than arbitrary software, so they will not be discussed in depth.

Finding OPT2: Attacker-authored VVM. Disabling O3 allows an attack in threat model A1 where the adversary injects a VVM that is verified by the director because the digital signature is not checked. Including O3, this attack also succeeds in threat model A*.

Finding OPT3: Supply chain attacks. We model supply chain attacks by placing a standard attacker on the system-level inputs (see Figure 2 attacker A#) and annotating every other connection as invulnerable to attack. With the supply chain attacker, all desired functional requirements are violated except for P2 (checks for VVM replays), which is maintained because there is no adversary between the primary ECU and director repository to replay VVMs. Supply chain attacks are out of scope for Uptane, but this modeling serves as a sanity check and speaks to the importance of supply chain security.